



# 游戏 编程精粹 2

GAME PROGRAMMING *Gems* 2



内附光盘



[美] Mark A. Deloura 编  
袁国忠 陈 萌 译  
姚 勇 审校



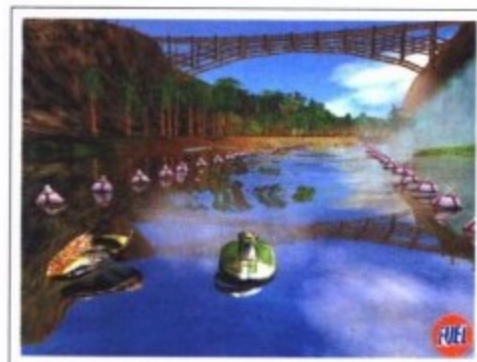
人民邮电出版社  
POSTS & TELECOM PRESS

# 游戏编程精粹 2

## GAME PROGRAMMING *Gems* 2

本书是“游戏编程精粹”系列丛书的第二本，包括70多篇全新的、探讨各种游戏编程主题的文章，读者可以将其中介绍的技术直接应用于自己的游戏项目中。这些文章全部由游戏编程专

家撰写，每篇文章要么提供了某个编程问题的实用解决方案，要么提出了一种创造性的减少编程时间和冗余的方法。该书由 *Game Developer* 杂志主编 Mark DeLoura 和各个游戏编程领域中出类拔萃的专家编辑，书中涵盖了开发最尖端的游戏引擎所涉及的主要主题。全书由内容丰富的六章组成。专家级开发人员可以立刻应用书中介绍的技巧，而初中级程序员通过阅读本书将增强技能，提高水平。本书是40多位经验丰富的游戏开发人员智慧和经验的结晶，是游戏开发必备的参考资料。



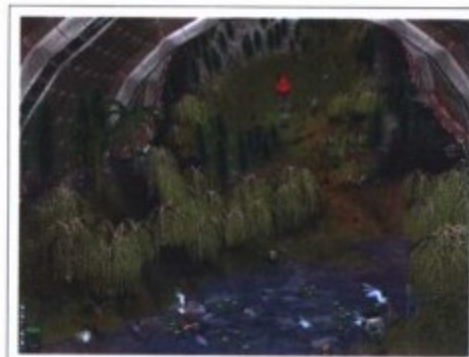
COPYRIGHT © 2001 FIRETOAD SOFTWARE



COPYRIGHT © 2001 PAST TREE GAMES



COPYRIGHT © 2001 GAS POWERED GAMES



COPYRIGHT © 2001 MUCKY FOOT PRODUCTIONS

ISBN 7-115-10871-4



9 787115 108715 >

人民邮电出版社网址 [www.ptpress.com.cn](http://www.ptpress.com.cn)

ISBN7-115-10871-4/TP·3190  
定价:75.00 元(附光盘)



---

## 内容提要

本书是“游戏编程精粹”系列丛书的第二本，包括 70 多篇全新的、探讨各种游戏编程主题的文章，每篇文章要么提供了某个编程问题的实用解决方案，要么提出了一种创造性的减少编程时间和冗余的方法。本书涵盖了开发最尖端的游戏引擎所涉及的主要主题，全书由 6 章组成，包括通用编程技术、数学技巧、人工智能、几何体管理、图形显示和音频编程。

专家级开发人员可以立刻应用书中介绍的技巧，而初中级程序员通过阅读本书将增强其技能和知识。这是一本必备的参考资料，是 40 多位经验丰富的游戏开发人员智慧和经验的结晶。

资源解密  
PDG

---

## 光盘内容介绍

本书附带光盘中包含书中列出的所有源代码、书中阐述的众多技巧的演示程序、DirectX 8.0a SDK、glSetup 单机版、OpenGL 实用程序工具包 (GLIT)、书中图像的高分辨率版本以及精彩的游戏开发网站的网址。用 C 和 C++ 编写的可移植的源代码概述性地阐述了大部分技巧，其中很多代码都是跨平台的，能够在 Macintosh 和 Linux 中正确运行，但有些技术只适用于 Windows/Xbox/DirectX 8.0a SDK。

### LIMITED WARRANTY AND DISCLAIMER OF LIABILITY

THE CD WHICH ACCOMPANIES THE BOOK MAY BE USED ON A SINGLE PC ONLY. THE LICENSE DOES NOT PERMIT THE USE ON A NETWORK (OF ANY KIND). YOU FURTHER AGREE THAT THIS LICENSE GRANTS PERMISSION TO USE THE PRODUCTS CONTAINED HEREIN, BUT DOES NOT GIVE YOU RIGHT OF OWNERSHIP TO ANY OF THE CONTENT OR PRODUCT CONTAINED ON THIS CD. USE OF THIRD PARTY SOFTWARE CONTAINED ON THIS CD IS LIMITED TO AND SUBJECT TO LICENSING TERMS FOR THE RESPECTIVE PRODUCTS.

CHARLES RIVER MEDIA, INC. ("CRM") AND/OR ANYONE WHO HAS BEEN INVOLVED IN THE WRITING, CREATION OR PRODUCTION OF THE ACCOMPANYING CODE ("THE SOFTWARE") OR THE THIRD PARTY PRODUCTS CONTAINED ON THE CD OR TEXTUAL MATERIAL IN THE BOOK, CANNOT AND DO NOT WARRANT THE PERFORMANCE OR RESULTS THAT MAY BE OBTAINED BY USING THE SOFTWARE OR CONTENTS OF THE BOOK. THE AUTHOR AND PUBLISHER HAVE USED THEIR BEST EFFORTS TO ENSURE THE ACCURACY AND FUNCTIONALITY OF THE TEXTUAL MATERIAL AND PROGRAMS CONTAINED HEREIN; WE HOWEVER, MAKE NO WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, REGARDING THE PERFORMANCE OF THESE PROGRAMS OR CONTENTS. THE SOFTWARE IS SOLD "AS IS " WITHOUT WARRANTY (EXCEPT FOR DEFECTIVE MATERIALS USED IN MANUFACTURING THE DISK OR DUE TO FAULTY WORKMANSHIP);

THE AUTHOR, THE PUBLISHER, DEVELOPERS OF THIRD PARTY SOFTWARE, AND ANYONE INVOLVED IN THE PRODUCTION AND MANUFACTURING OF THIS WORK SHALL NOT BE LIABLE FOR DAMAGES OF ANY KIND ARISING OUT OF THE USE OF (OR THE INABILITY TO USE) THE PROGRAMS, SOURCE CODE, OR TEXTUAL MATERIAL CONTAINED IN THIS PUBLICATION. THIS INCLUDES , BUT IS NOT LIMITED TO, LOSS OF REVENUE OR PROFIT, OR OTHER INCIDENTAL OR CONSEQUENTIAL DAMAGES ARISING OUT OF THE USE OF THE PRODUCT.

THE SOLE REMEDY IN THE EVENT OF A CLAIM OF ANY KIND IS EXPRESSLY LIMITED TO REPLACEMENT OF THE BOOK AND/OR CD-ROM, AND ONLY AT THE DISCRETION OF CRM.

THE USE OF "IMPLIED WARRANTY" AND CERTAIN "EXCLUSIONS" VARY FROM STATE TO STATE, AND MAY NOT APPLY TO THE PURCHASER OF THIS PRODUCT.



---

## 关于翻译审校

姚勇，H3D 游戏开发公司创始人。小学四年级因痴迷游戏，而被父母教导学习编程就可以自己写游戏玩（成年人的花招），姚勇便从 BASIC 入手，开始学习自己编写游戏。为了提高游戏速度，他六年级自学 Z80 芯片汇编语言，初中自学 APPLEII 汇编。在桌面 PC 电脑 3D 图形化的革命风暴来临之季，3DFX Voodoo 3D 加速卡进入中国、Direct3D 版本 5 产生、quake1 出现，姚勇自己编制游戏的渴望又再次被强烈唤起。1998 年他从清华大学毕业并留校任教，工作之余，为了其希望在中国制作出世界级三维游戏的理想，他开始筹备 Hardcore3D（H3D）游戏开发公司。2000 年 8 月 Hardcore3D Studio 正式成立，吸引了许多国内三维引擎精英，姚勇也辞去清华大学的工作，专心投入到游戏开发事业中来。在经历 3 年努力后，姚勇终于带领着 H3D 摆脱了实时渲染技术以及超大多用户网络服务器分布技术困扰，开始制作 3D MMORPG 引擎。目前他们正在制作一款 3D MMORPG 游戏，希望能真正改变中国游戏制作落后的现状，形成中国游戏制作技术之硬核。

H3D 公司的网址：<http://www.hardcore3d.net>。

姚勇的 Email 地址：[puzzy@hardcore3d.net](mailto:puzzy@hardcore3d.net)。



---

# 前言

Mark DeLoura  
madsax@satori.org

欢迎阅读《游戏编程精粹 2》。本书精选了 70 篇讨论各种游戏编程主题的文章，读者可将其中阐述的技术直接应用到游戏中。能给读者呈现一本这样的图书，我感到自豪。

开始策划该书时，我的目标就是确保其内容尽可能地精彩而实用。我不希望本书是一本充斥着理论或过时信息的图书，因此，在组织本书时，我没有闭门造车，而是邀请专家担任各章的责任编辑，由他们负责挑选探讨相关主题的文章。所幸的是，给《游戏编程精粹 1》撰过稿的很多开发人员都有意负责本书某章的组稿工作。

下面是本书各章的责任编辑。他们工作繁忙，但为确保每篇文章都值得阅读付出了大量的心血。

- 通用编程技术：Gas Powerd Games 公司的 Scott Bilas;
- 数学技巧：Naughty Dog 公司的 Eddie Edwards;
- 人工智能：任天堂（美国）公司的 Steve Rabin;
- 几何体管理：C4 Engine 公司的 Eric Lengyel;
- 图形显示：Nvidia 公司 D. Sim Dietrich Jr.;
- 音频编程：Lithtech 公司的 James Boer。

## 幕后故事

---

组织《游戏编程精粹 1》时，我还是任天堂（美国）公司的主力软件工程师。我很灰心，游戏行业的公司都在制定封闭的标准、创建封闭的 API、申请软件专利，在自己和他人之间砌起一道无形的墙。坦率地说，我对整个行业极其失望。在给出版社发送《游戏编程精粹 1》稿件的前一天的凌晨 3 点，我再次阅读了刚完稿的“前言”。我在其中写道，制定开放标准并同行业中的其他人共享知识非常重要。就在这一刻，我意识到了自己希望致力于促进知识在游戏业中的传播，而不是宣传另一个专用标准的强烈愿望。

不久后，经过同 E3 的一次独特合作和多次深入交流后，我决定接受 Game Developer 杂志的主编职务。现在，我每月都给读者呈上几篇精彩的文章！



## 《游戏编程精粹 2》的构思

---

正如前面指出的，本书是本人同 6 位责任编辑合作的结晶。我们旨在确保每章的内容对专家级游戏开发人员来说都颇有参考价值。如果读者是初中级游戏开发人员，请仔细阅读书中的文章，学习其中的技巧；如果需要更详细了解某个主题，可阅读相应的参考文献。书中的很多文章只对复杂的技巧做了概要性的介绍，但提供了参考文献，读者可通过这些参考文献更深入地了解相应的主题。当然，网上有很多不错的游戏开发资源，其中包括 gamasutra.com、flipcode.com 和 gamedev.net。

我需要强调的一点是本书中 API 的使用。在《游戏编程精粹 1》中，所有的文章使用的都是 C/C++ 和 OpenGL，旨在确保代码能够在当前的各种系统上运行。本书后退了一步。由于当前的游戏平台非常多（包括游戏控制台、计算机、手持设备、Web 和无线设备），要编写能够在所有平台上运行的代码几乎是项无法完成的任务。然而，对于任何称职的开发人员来说，有些语言和 API 是必须熟知的，这包括 C、C++、x86 汇编语言（如果您是 PC 引擎程序员）、OpenGL 和 Direct3D（如果您是 3D 程序员）。您可能还应该精通 Java 和一些脚本编写语言。我们假设作为专家级游戏开发人员的读者在理解这些语言时不会有任何困难，书中的有些地方可能还会超过以上的知识范围。

## 游戏行业

---

多年来，游戏开发行业一直呈稳定增长的态势。我们喜欢引用的一项统计数据是：在美国，游戏行业的收入已超过了电影业。这样的成绩好像是令人高兴的，我们应暂时放下手中的编程工作，为此狂欢。但是别忘了，电影票的平均价格是 8 美元，而每个游戏软件差不多要 40 美元！因此，实际上游戏玩家的人次只有电影观众的五分之一。

游戏行业的增长空间还很大。将最新的 *Miyamoto* 游戏下载到宽带机顶盒中并同地球另一边的朋友一起玩的日子还未到来。然而，这指日可待。也许是几年后，这将打开通向各种新娱乐方式的大门。

世界各地的人都在玩游戏和做游戏。我收到的来信中，有的来自俄罗斯、克罗地亚、澳大利亚、新西兰、巴西、欧洲和北美地区的开发人员。游戏已成为一种全球景象，人们通过共同的体验走到一起。当您设计下一个游戏时，不仅应考虑跨平台的兼容性，还应考虑跨文化的兼容性！

随着游戏在全球的影响越来越大，我们需要注意的一个问题是，由于游戏玩家已非常多，游戏已引起了文化监察人员的注意。每当“问题少年”实施暴力犯罪时，我们都“荣幸”地位列在某种程度上必须对此负责的媒体之一：电影、电视、音乐和游戏。这虽然令人不快，但它表明游戏已成为公众的话题之一：游戏不再只是供小孩玩的东西，而是一种真正的艺术形式，因此同其他艺术形式一样，游戏也是被监督的对象。

媒体中的暴力是否会增强人们的暴力倾向这一问题也许永远不会有答案，但在公众看来，答案是肯定的。作为一个对此负有责任的行业，我们力所能及的重要工作之一就是提倡公众接受 ESRB 分级制度，让父母能够选择孩子该玩什么样的游戏，进而担负起他们应付的

责任。

当然，通过玩游戏确实能够学到一些知识；否则，玩家如何能够通过第 20 级中的怪物把守的关口？玩家学到的知识随游戏而异。随着游戏与现实世界越来越接近，玩家在游戏中受教育的可能性也越来越大。您将教给玩家什么样的知识？

## 开放标准

---

当前，语言和 API 越来越多，其中很多是专用的，这是最令我沮丧的事情之一。独立开发人员在选择平台时面临的挑战日益严峻，因为难以获得准确的有关开发软件包的信息。游戏行业中对开发系统保密的潮流让我难以理解。

在 PC 上，这不是什么大问题。在 PC 上面，系统正越来越开放。然而，令人惊讶的是，大多数技术创新是在 PC 上进行的。几个最新的游戏在发布后便立刻公开了源代码。

几年前，哈佛商学院对施乐 PARC 的研究人员进行了一项调查。调查旨在找到确保研究人员成功的最重要的品质，结果发现交流和共享信息的能力是确保成功最重要的品质之一。离群索居的研究人员可能干得不错，但真正卓越的研究人员都有丰富的专业网，经常同其他人员交流。

上述调查结果也适用于游戏行业。当前，一些最成功的游戏开发公司都公开其游戏，同每一个人交流。他们公开游戏源代码、关卡编辑工具和 API，他们在游戏开发人员会议中发表演讲。他们的游戏玩家得以创建如梦如幻的新关卡和游戏风格 (mod)，而求知若渴的游戏开发人员得以向专家学习。通常，人们认为公开意味着失去，但在游戏行业，将您的知识产权公开，也许是您能做出的最为明智的决策之一。如果能让游戏玩家成为游戏开发人员，无疑是个多赢的结局。





---

## 致 谢

参与本书制作的人员为数众多，这里只能感谢其中的几位。首先（也是最重要的），要感谢 Charles River Media 出版社的 Dava Pallai 和 Jenifer Niles，在他们的鼓励下，我才继续策划“游戏编程精粹”系列丛书。本书的面世得益于他们对卓越品质的不懈追求和始终一丝不苟的要求。感谢 Andrew Glassner 提出了最初的构思。*Graphics Gems* 系列丛书现仍在我的书架上占据一席之地，如果读者还未购买，请赶快去。

真正创作本书的是各章的责任编辑和撰写文章的作者们。他们大部分都有专职工作，很多还工作繁忙，但为确保本书出类拔萃，他们放弃了很多休息时间。这里衷心感谢他们将自己的专业知识拿出来与大家共享。

最后，要感谢的第三批人是我在 Gama Network 的同事们。与他们共事绝对是令人惊喜的。我每天都同这些全心全意地致力于使游戏开发行业繁荣的人一起工作。每天同像 Jennifer Pahlka、Alan Yu、Susan Marshall、Alex Dunne 和 Dan Huebner 这样的人一起工作让人感到快乐。他们确保游戏开发人员会议、Gamasutra.com 和 *Game Developer* 杂志平稳地运转，并提供丰富而精彩的信息。这里还要感谢 *Game Developer* 杂志的全体编辑人员，他们是 Laura Huber、Jennifer Olsen 和 Tiger Byrd。他们积极的工作态度和优秀的职业素养让我惊讶。感谢他们帮助一名工程师学会了如何成为一名成功的主编！



---

## 关于封面图案

Andrew Kirmse



封面上的图片是在引擎内部截取的飞行动作游戏《星球大战之星际战士》（*Star Wars Starfighter*）中的一个场景，这是 LucasArts 于 2000 年 2 月发布的一款 Playstation 2 游戏。其中的地形、岛屿和瀑布是使用 3D Studio Max 创建的，建筑物和敌军是由关卡编辑工具（level editing tool）加入的。飞行器使用二次纹理（second texture）实现倒影；地形使用二次细节纹理（second detail texture）实现。瀑布下的薄雾、爆炸和烟雾以及船艇的尾迹和激光是使用粒子系统模拟的。

在这个场景中，玩家的轰炸机（Nym 的中型轰炸机）摧毁了一架敌机，其中的精华是充斥碎片的画面。敌机能够同时绕过地面障碍物、瞄准目标、避开飞船和追击玩家——但要命的是，它躲不开 Nym 的激光。

星河网  
PDG



---

## 作者简介

### Scott Bilas

---

scottb@aa.net

Scott Bilas 是 Gas Powered Games 公司的一名编码人员, 享受着编写后端游戏系统的乐趣。他偶尔还从事写作, 曾在 *Game Programming Gems* 和 *Game Developer* 杂志上发表过文章。Scott 最近承担了《游戏编程精粹 2》第 1 章的编辑工作, 撰写了书中的几篇文章, 在 GDC 2001 上发表演讲, 开发 *Dungeon Siege*。

### James Boer

---

jimb@lithtech.com

James Boer 是 Lithtech Inc. 的一名音频程序员, 参与了 *Deer Hunter I & II*、*Rocky Mountain Trophy Hunter*、*Microsoft Baseball 2000* 和 *Tex Atomic's Big Bot Battles* 等项目的开发工作。

### Charles Cafrelli

---

skywise@iquest.net

Charles 毕业于普渡大学, 在 Origin Systems (R.I.P.) 公司从事过多个游戏的开发工作, 其中包括 *Wing Commander 2*、*Privateer*、*Wing Commander 3* 和 *Pacific Strike*。他还在 Boehringer-Mannheim 开发过用于血糖仪的软件, 当前正在开发用于下一代消费电子设备的软件。他喜欢的游戏包括 Scott Adams 的 *Adventure*、Steve Meretsky 的 *Planetfall*、*Elite*、*Space Ace*、*Command and Conquer*、任何有 Zelda 的游戏以及 *Final Fantasy*。

### Simon Carter

---

scarter@bigbluebox.com

Simon Carter 1986 年就进入了计算机游戏行业, 那时他才 11 岁, 帮助哥哥 Dene 开发 Sinclair ZX Spectrum。1994 年, 他受聘于 Bullfrog Productions Ltd, 在开发 *Magic Carpe* 的过程中进一步提高了编程技能, 继而成为

*Dungeon Keeper* 的项目负责人和主力程序员。1999 年, Simon 与兄弟一起创建了自己的游戏公司 Big Blue Box Studios, 现正从事一个还未发布的运行在 Microsoft Xbox 中的游戏项目。

## Peter Dalton

---

pdalton@xmission.com

www.xmission.com/~pdalton

Peter Dalton 现为 Evans & Sutherland 的高级显示小组的一员, 该小组致力于开发下一代视觉显示 (visual display)。他的工作重点是提高用于训练飞行员的飞行仿真装置的质量。

## Bruce Dawson

---

bruced@humongous.com

www.cygnus-software.com/papers/

Bruce Dawson 是 Humongous Entertainment 公司的技术总监, 其工作是完成其他人没有时间做的各种有趣而颇具挑战性的任务。业余时间, Dawson 还在 DigiPen 从事培训工作。加入 Humongous Entertainment 公司之前, Dawson 受聘于 Cavedog Entertainment 公司, 协助各种产品小组完成工作。Dawson 在 Avid Technology 公司的 Elastic Reality 分部工作过多年, 从事特效软件和视频编辑插件的开发。很久以前, 他还在 Electronic Arts (加拿大) 公司工作过, 那时该公司名为 Distinctive Software。

## Mark DeLoura

---

maidsax@satori.org

www.satori.org/maidsax

Mark 是 *Game Developer* 杂志的主编, 也是《游戏编程精粹 1》和《游戏编程精粹 2》的编辑。加入 *Game Developer* 之前, Mark 在任天堂 (美国) 公司工作过 5 年, 担任开发人员关系小组的主力软件工程师, 从事 Gamecube 和 Nintendo 64 开发。在此之前, Mark 还从事过虚拟现实方面的研究, 是多个 Usenet 虚拟现实新闻组的主持人之一。

## D. Sim Dietrich Jr.

---

sdietrich@nvidia.com

D. Sim Dietrich Jr. 从 Apple II+ 开始, 就从事实时图形显示的研究, 现为 Nvidia 公司技术开发人员关系小组的一员, 帮助 PC 开发人员开发尖端的实时 3D 图形。他当前感兴趣的课题包括逐像素光照技术 (per-pixel lighting) 和阴影 (shadows) 以及使用软件工程技术来提高游戏开发过程的生产率和创造性。

## Nathan d'Obrenan

---

nathand@firetoads.com

www.firetoads.com/coding

Nathan 在小学时涉足视频游戏编程，完全靠自学成才。Nathan 是位于加拿大卡尔加里市的游戏公司 Firetoad Software, Inc 的主力程序员。其工作包括研究和开发全新的技术，以便将其用于游戏中。

## Carl Dougan

---

cdougan@adrenalin.com

Carl Dougan 有 6 年从事游戏开发的经验，现在加利福尼亚从事 XBox 游戏 MassMedia 的开发工作。

## Eric Dybsand

---

eric@geta.cncdsl.com

Eric Dybsand 是《游戏编程精粹 1》的撰稿人之一，1987 年便进入了游戏开发行业。他为 2001 年发布的多款赛车游戏开发了赛车 AI。在 2000 年，他还为太空战略游戏 *Master of Orion 3* 提供战略 AI 方面的咨询。Eric 还为棒球和摔跤游戏提供过 AI 方面的咨询。他为 Windward Studios 的 RTS 游戏 *Enemy Nations* 以及 Fenris Wolfe 的 FPS 游戏 *Rebel Moon Revolution* 和 *War In Heaven* 设计和开发了 AI 对手 (opponent)。在过去的前五届游戏开发人员会议中，Eric 担当了 AI 圆桌会议的主持人和 CGDC 演讲人；他当前是 Gamasutra 的“AI 编程”论坛和 IGDA 的“游戏中的 AI”论坛的主持人。

## Eddie Edwards

---

eddie@tinyted.net

Eddie Edwards 进入游戏编程行业之前，在剑桥大学学习数学，并做过一段时间的软件工程师。1999 年加入位于加利福尼亚的 Naughty Dog 公司之前，他在英国从事过各种游戏项目，包括 *Wolfenstein 3D* 和 *DOOM*。他当前是将于 2001 年圣诞节发布的 Playstation 2 游戏的编程小组成员。

## Thomas Engel

---

thomas.engel@factor5.com

Thomas Engel 是 Factor 5 公司的技术总监，也是该公司的创始人之一。他有 10 多年游戏开发的从业经验，研究重点是“玩游戏”后面的技术。在过去的一年中，他负责的项目之一



是 Factor 5 的 *MusyX* 声音系统。Thomas 当前正忙于开发用于下一代控制台（如任天堂的 Gamecube）的音频和图形技术。

## Jeff Evertt

---

jeff@evertt.com

Jeff 在 Intel 从事过多年软件和硬件方面的设计，还在 Cavedog Entertainment 公司从事过多个游戏项目，现在 Littech Inc 开发 PC 和 Xbox 渲染程序（renderer）。

## Mark Fischer

---

beach@beachsoftware.com

Mark Fischer 创建了开发加密软件的 Beach Software 公司。Mark 还是一名独立的咨询人员，擅长优化网络软件和确定性软件架构。

## Tom Forsyth

---

tomf@muckyfoot.com

在其 ZX Spectrum 上看过 *Elite* 后，Tom 就被 3D 图形迷住了，他一直致力于榨干硬件的最后一滴油。他在 Spectrum、Sinclair QL、Atari ST、Sega 32X、Saturn、Dreamcast 和 PC 上编写过绘制三角形的函数，而现在已感到厌倦。他非常感激我们有能够为他绘制图形的硬件。在 3Dlabs 从事两年编写 3D 图形卡驱动程序的工作后，Tom 加入了 Mucky Foot Productions——一家位于英国吉尔福德市的游戏公司。

## Dan Ginsburg

---

ginsburg@alum.wpi.edu

Dan Ginsburg 现为 ATI Research 公司的 OpenGL 驱动程序小组的一名软件工程师，主要负责实现扩展，创建演示程序和编写新的扩展规范。受聘于 ATI 之前，Dan 在 n-Space 公司待了两年，从事 Playstation 和 PC 游戏 *Die Hard Trilogy 2: Viva Las Vegas* 的开发，是引擎小组的成员之一。

## Miguel Gomez

---

miguel@littech.com

Miguel Gomez 是交互式应用程序软件开发商 Littech Inc. 的首席工程师，专门从事数值分析和物理建模。他以前的工作经历包括图形编程和物理学编程，从事的项目有 Electronic Arts 的 *PGA Tour '96*、Activision 的 *Hyperblade*、Microsoft 的 *Baseball 3D* 和 Psygnosis 的 *Destruction Derby 64*。他获得了华盛顿大学的物理学学士学位，现正在攻读该大学的应用数

学硕士学位。

---

## Dave Gosselin

---

[gosselin@ati.com](mailto:gosselin@ati.com)

Dave Gosselin 现为 ATI Research 的 3D 应用程序研究小组的软件工程师，参与各种演示程序和 SDK 工作，编写 OpenGL 扩展规范。在此之前，他在多家公司工作过，其中包括 Oracle、Spacetec IMC、Xyplex 和 MIT 林肯实验室，从事过从底层网络技术和 Web 技术到图形处理和 3D 输入设备等各种项目。

---

## Bryon Hapgood

---

[bhapgood@kodiakgames.com](mailto:bhapgood@kodiakgames.com)

Bryon 在 Bullfrog Productions 开始其职业生涯，从事用于 Amiga 的 *Populous II*；现为 Kodiak Interactive 的一名程序员。

---

## Evan Hart

---

[ehart@ati.com](mailto:ehart@ati.com)

Evan Hart 是 ATI 的 3D 应用程序研究小组的一名软件工程师，当前的研究重点是 OpenGL 软件和扩展。加入 ATI 之前，他在 Battelle 从事基于 PC 的 3D 模拟。

---

## Pete Isensee

---

[pkisensee@msn.com](mailto:pkisensee@msn.com)

[www.tantalon.com/pete.htm](http://www.tantalon.com/pete.htm)

Pete 编写过各种游戏：从 7-CD 游戏、多玩家游戏到控制台应用程序。他拥有计算机工程学位，并忙里偷闲地研究 C++ 优化技术。

---

## John Isidoro

---

[jisidoro@cs.bu.edu](mailto:jisidoro@cs.bu.edu)

John Isidoro 是 ATI Technologies 的应用程序研究小组的 3D 图形编程顾问，正在攻读波士顿大学的计算机科学博士学位。需要指出的是，John 确实喜欢四元数。

---

## Greg James

---

[gjames@nvidia.com](mailto:gjames@nvidia.com)

Greg James 是 Nvidia 的技术开发人员关系小组的软件工程师，从事实时 3D 图形工具和

演示程序的开发。在此之前，他是一家高能物理实验室的研究助理。直到有一天晚上，他觉得自己对不起被裂变的质子，因此决定离开，去从事 PC 游戏和计算机图形方面的工作——自从老爸将一个名为 Amiga 1000 的米色怪物带回家，这便是他学生时代的一项业余爱好。

## Lasse Staff Jensen

---

lasse@funcom.com

www.funcom.com

Lasse Staff Jensen 获得了挪威的特隆赫姆工程学院的计算机工程学位。4 年前进入游戏行业之前，他是陆军工程学院的助教兼系统管理员。Lasse 在计算机图形学方面的造诣颇深，拥有将近 10 年的图形应用程序编程经验；他现为 Funcom 的技术经理，带领程序员小组研究和开发 Funcom 的核心技术。

## Yossarian King

---

yking@ea.com

Yossarian King 是 Electronic Arts（加拿大）公司的一名工程师，是 *FIFA Soccer* PC 版的主力程序员；之前还从事过 Nintendo 64、3DO、Sega Saturn 和 Commodore 64 方面的开发。

## Andrew Kirmse

---

ark@alum.mit.edu

Andrew 是 *Meridian 59*（1996）的导演和编剧之一，还在 *Star Wars: Starfighter*（2001）中担任图形程序员。他拥有 MIT 的物理、数学和计算机科学学位。Andrew 现受聘于 LucasArts Entertainment 公司。

## Jesse Laeuchli

---

jesse@laeuchli.com

Jesse Laeuchli 是一名获得认证的 C/C++ 和 Windows 程序员，现在匈牙利的布达佩斯定居。身为外交官之子，他在很多国家和地区居住过，如中国、沙特阿拉伯和非洲的一些国家。他给多本计算机杂志和网站撰稿。在计算机编程方面，他感兴趣的领域是计算机图形和操作系统。他现在匈牙利学习，同时从事 OpenGL 3D 引擎的开发。他还管理了一个小型网络，从事一般性的计算机故障诊断和排除工作。

## Adam Lake

---

adam.t.lake@intel.com

Adam Lake 是 Intel Architecture Lab (IAL) 的图形和 3D 技术小组的一名高级软件工程师。

他获得了北卡罗莱纳大学的硕士学位。

## François Dominic Laramée

---

francoislaramee@videotron.ca

pages.infinet.net/idjy

François Dominic Laramée 有近 10 年的游戏从业经验，在此期间，他独立开发了 20 多个运行在控制台、PC、网络、电视机顶盒和各种怪异平台上的游戏。他撰写的关于游戏设计、编程和生产的文章总能吸引成千上万的读者；他还“骗取”两所不同的大学授予他研究生学位（确切地说是一所，另一个学位还在攻读中）。

## Eric Lengyel

---

lengyel@c4engine.com

Eric Lengyel 是 *C4 Engine* 项目的主力程序员，还是 *Mathematics for 3D Game Programming & Computer Graphics* (Charles River Media, 2001) 一书的作者。之前，Eric 曾是 Sierra 的 3D 引擎设计师和 Apple Computer 的 OpenGL 工程师。他拥有弗吉尼亚理工学院的数学硕士学位。

## Ian Lewis

---

ilewis@acclaim.com

Ian Lewis 在 1993 年就进入了计算机娱乐行业，现为位于犹他州盐湖城的 Acclaim Studios 的一名网络和音频程序员。

## Noel Llopis

---

nllopis@mgigames.com

www.cs.unc.edu/~llopis/

Noel Llopis 在 1985 年便开始在 Amstrad CPC 上编写游戏，将其作为一项业余爱好。他获得了马萨诸塞州阿默斯特大学的计算机工程学士学位和北卡罗莱拉大学的计算机科学硕士学位。几年前，他终止博士学业，进入游戏行业；现为 Meyer/Glass Interactive 的一个 Xbox 游戏项目的主力工程师。

## John Manslow

---

jfm96r@ecs.soton.ac.uk

John Manslow 于 1994 年以优异的成绩在英国诺森伯兰大学获得微电子工程学位，专业为硬件人工神经网络。在该大学的通信研究小组工作了短暂的时间后，他加入了 Neural Computer Sciences，为开发灵活的面向对象神经网络库发挥了重要作用。1996 年，John 进



入南安普顿大学的图像、语音和智能系统研究小组攻读博士学位，研究方向为使用 AI 技术分析卫星图像。完成博士学业后，他加入了 Codemasters Software Company Limited，担当一名 AI 研究和开发程序员，从事了多个旨在将游戏学习技术引入到 Codemasters 未来的产品中的项目。2001 年初，John 跳槽到 Neural Technologies Limited，从事使用高级 AI 技术检测电信诈骗（telecommunications fraud）的研究。当前，他感兴趣的课题包括 AI 在游戏中的应用和人工组织的进化。

## Herb Marselas

---

hmarselas@ensemblestudios.com

www.ensemblestudios.com

Herb Marselas 是 Ensemble Studios 的一名 3D 引擎专家，致力于下一代即时策略游戏的开发。

## Carl S. Marshall

---

carl.s.marshall@intel.com

Carl S. Marshall 是 Intel 架构实验室（Architecture Labs）图形和 3D 技术小组的一员，拥有 Clemson 大学计算机科学硕士学位，在攻读硕士期间，从事虚拟现实的研究。他花了一个暑假在领先的游戏开发公司 Valve 了解游戏行业的方方面面。当前，Carl 正研究图形引擎 Shockwave3D 的各个方面，同时致力于人工智能和创建实时 3D 非照片真实感（non-photorealistic）渲染的虚拟世界和效果。

## Chris Maughan

---

cmaughan@nvidia.com

Chris 有 6 年 3D 图形硬件行业的从业经验。最初，他主要从事设备驱动程序的开发，加入 Nvidia 开发人员关系组之前，他在 3Dlabs 领导 DirectX 驱动程序小组。从驱动程序领域转向图形管道（graphics pipeline）领域给他提供了致力于软件开发的机会，这样每当出现 bug 时，无须重新启动机器。在 Nvidia 开发人员关系小组，Chris 专门研究协助游戏开发的工具。他当前的工作重点是调试和生成用于 DirectX 的着色器。

## Jason L. Mitchell

---

jasonm@ati.com

Jason L. Mitchell 是 ATI Research 的 3D 应用程序研究小组的负责人，从 1996 年起，他就该小组工作。3D 应用程序小组从事的工作位于硬件、应用程序和 API 开发的交叉领域，需要同 OpenGL ARB、Microsoft Direct3D 小组和 ATI 内部的硬件设计师交流。该小组还开发新颖的渲染技术和 ATI 的所有技术演示程序。加入 ATI 之前，Jason 在辛辛那提大学从事人

眼跟踪技术在人机界面中的应用研究，并获得电子工程硕士学位。他于 1994 年获得 Case Western Reserve 大学的计算机工程学士学位。除“游戏编程精粹”系列外，Jason 还在 *Game Developer* 杂志和 Gamasutra.com 上发表过文章。

## Aaron Nicholls

---

aaron\_feedback@hotmail.com

Aaron Nicholls 是位于华盛顿州雷德蒙市的 Microsoft 公司的开发小组的负责人。他感兴趣的领域是计算机图形学和 AI，并一直在这些领域寻找新的挑战。他欢迎您的反馈和来信。

## John Olsen

---

infix@xmission.com

John Olsen 于 1989 年获得犹他大学计算机科学系授予的学士学位，1997 年进入游戏开发领域之前，从事图形库和演示软件的开发。他参与了包括 *Jet Moto 2* 在内的多个 Playstation 游戏的开发，现正帮助 Microsoft 公司开发一个 Xbox 游戏。他感兴趣的领域包括自主 AI 行为、网络技术、数据组织和跟踪，并坚定不移地讨伐特例代码（special-case code）。

## Javier F. Otaegui

---

javier@sabarasa.com

www.sabarasa.com

Javier F. Otaegui 是位于布宜诺斯艾利斯的游戏开发公司 Sabarasa Entertainment 的 CTO。他还是 RTS 游戏 *Malvinas 2032* 的设计师和主力程序员，该游戏已经在阿根廷和美国销售，不久后将在欧洲销售。Javier 当前感兴趣的领域包括 3D LOD 风景、网络多玩家技术、让人成瘾的“玩游戏”以及物理系统和模拟。要获得更详细的有关 *Malvinas 2032* 的信息，请访问该公司的网站或给 Javier 发 E-mail。

## Kim Pallister

---

kim.pallister@intel.com

Kim Pallister 是 Intel 解决方案实现（solutions enabling）小组的技术营销工程师和处理器宣传员（evangelist），当前致力于实时 3D 图形技术和游戏开发。欢迎读者的 E-mail。

## Scott Patterson

---

scottp@tonebyte.com

过去的 10 年中，Scott Patterson 受聘于 Naughty Dog、Midway 和 Microprose。在此期间，他从事的领域包括音频系统、图形代码和开发工具；现正领导游戏引擎的开发。

## Matt Pritchard

---

mpritchard@ensemblestudios.com

Matt Pritchard 是 Ensemble Studios 的一名开发人员，是开发销售量达数百万的 PC 游戏《帝国时代》的程序员之一，他一直深入参与了该系列游戏的开发。另外，他还撰写了大量有关游戏开发的文章，涉及的主题包括优化、图形、多玩家“作弊”技术（multiplayer cheating）等。也许是《游戏编程精粹 1》的内容过于简单，他将文章的交稿日期和第一个孩子的预产期安排在同一周。不编写游戏时，他同家人呆在一起（家庭成员包括女儿和一条黑色的拉布拉多猎狗），或纵情于业余爱好，如收集过时的计算机和视频游戏。如果您需要 Atari 2600 *Quadrant* 游戏、Exidy Sorcerer 计算机或对他撰写的文章有何看法，请给他发 E-mail。

## Steve Rabin

---

steve\_rabin@hotmail.com

Steve Rabin 是《游戏编程精粹 1》的撰稿人，现为《游戏编程精粹 2》人工智能一章的责任编辑和撰稿人。Steve 现为任天堂（美国）公司 Gamecube 软件开发支持小组的成员，从事演示程序编写、新技术研究和帮助其他开发人员的工作。加入任天堂之前，Steve 在多个小型的新兴公司工作过，是 3 款游戏的 AI 工程师。他曾在游戏开发人员会议上就 AI 发表过演讲，有 10 年视频游戏行业的从业经验。他获得了华盛顿大学计算机工程学士学位，专业为机器人技术。

## John W. Ratcliff

---

jratcliff@verant.com

John W. Ratcliff 有 20 多年软件开发和写作经验。John 参与开发的产品包括 Electronic Arts 的 688 *Attack Sub*、SSN-21 *Seawolf* 和 *Scarab* 以及 989 Studios 的 *Cyberstrike 2*。除上述游戏外，John 还为 Milliken Publishing Company 和 Cardiosoft 分别开发过众多的教育软件和生物医学软件系统。通过开发 *Digpak*、*Midpak* 和 *Vidpak*，John 为数百个游戏的开发做出了技术方面的贡献。John 现为 Sony Online Entertainment 的高级技术设计师，从事多玩家第一人称射击游戏 *Planetside* 的开发。

## Graham Rhodes

---

rhodes@sed.ara.com

Graham Rhodes 是位于新墨西哥州阿尔伯克基市的工程研发公司 Applied Research Associates, Inc. (ARA) 的 Quicksand Games 分公司的一名高级科学家。Graham 有 10 多年计算机几何学、交互和实时 3D 图形学以及物理学编程方面的经验。在 ARA 期间，他参与了各种项目，其中包括游戏中的实时计算机流体动力学和飞碟的非线性仿真。从 1996 年到 2000

年, Graham 是 NASA 的下一代革命性分析和设计环境(Next Generation Revolutionary Analysis and Design Environment, NextGRADE) 图形用户界面的签约主力开发人员和软件设计师, 负责实现计算机图形和拓扑算法。

## Jason Shankel

---

shankel@pobox.com

Jason Shankel 是 Maxis/EA 的一名游戏程序员。除《游戏编程精粹 1》外, 他还给 *Dr. Dobbs Journal* 和 *DevJournal* 撰稿。

## Greg Snook

---

gregsnook@home.com

Greg 在视频游戏行业摸打滚爬了 6 年多。他最初从事动画工作, 但由于被编码错误困扰而转向编程领域。从此以后, Greg 在多款 PC 和控制台游戏项目中担当了美工和程序员。

## Bruno Sousa

---

gems2@fireworks-interactive.com

www.fireworks-interactive.com

Bruno 有超过 8 年的编程经验, 在 17 岁时就获得了第一份游戏编程工作。他现正在创建自己的公司, 致力于为游戏开发公司开发廉价的工具和游戏预算 (budget games)。

## Marco Tombesi

---

tombesi@infinito.it

digilander.iol.it/baggior/

Marco Tombesi 是一位刚毕业的意大利工程师, 擅长 C++、OpenGL、Java 和 CORBA, 管理过一些自由游戏项目, 积累了丰富的经验 (但两手空空)。他喜欢游戏, 尤其是体育方面的, 当前正致力于开发一款新的足球游戏, 不久后就可能推出, 到时任何人都能够享受他的工作成果。

## Paul Tozour

---

gehn29@yahoo.com

Paul Tozour 当前是 *Thief 3* 项目的一名 AI 开发人员。他最新的工作成果之一是为广受欢迎的游戏 *MechWarrior 4: Vengeance* 开发搏击 AI。



## William van der Sterren

---

[william@cgf-ai.com](mailto:william@cgf-ai.com)

William van der Sterren 是 CGF-AI 的老板和 AI 开发人员。他研究和开发用于游戏和仿真的战术 AI, 并在游戏开发人员会议上发表过演讲。William 现为飞利浦电子研究实验室(Philips Electronics Research Labs) 的一名设计师/科学家, 从事 AI 方面的工作。

## Alex Vlachos

---

[alex@vlachos.com](mailto:alex@vlachos.com)

[alex.vlachos.com](http://alex.vlachos.com)

Alex Vlachos 现为 ATI Research 的 3D 应用程序研究小组成员, 从 1998 起, 他便在这里工作。Alex 是 ATI Radeon 图形处理器发布的演示程序 *Radeon's Ark* 的主力开发人员, 并继续编写展示下一代硬件特性的 3D 引擎。另外, 他还开发了 N-Patches (一种曲面表示法, 是 Microsoft DirectX 8 的一部分)。加入 ATI 之前, 他在 Spacetec IMC 工作, 是 SpaceOrb 360——一个自由度为 6 游戏控制器——的软件工程师。Alex 毕业于波士顿大学。

## Scott Wakeling

---

[scott@chronicreality.com](mailto:scott@chronicreality.com)

[www.chronicreality.com](http://www.chronicreality.com)

Scott 是 Virgin Interactive 公司的一名程序员, 还是独立的游戏小组 Chronic Reality 的创始人, 该小组致力于探索一切有趣的东西, 而不管它从技术上说有用还是没用。如果您想让他迷途知返, 请给他发 E-mail。

## Keith Weiner

---

[keith@dw.com](mailto:keith@dw.com)

Keith Weiner 是 DiamondWare 公司的 CEO。从 1994 年起, 他设计过音频产品、引擎、管道 (pipeline) 和设备驱动程序。

## Steven Woodcock

---

[ferretman@gameai.com](mailto:ferretman@gameai.com)

[www.gameai.com](http://www.gameai.com)

Steven Woodcock 的游戏 AI 经验源自 18 年的弹道导弹防御工作, 在那期间他开发了大量实时战争游戏和模拟程序。在民用领域从事短暂的工作后, 他重返国防领域, 帮助国家导弹防御系统开发 AI。Steve 的网站致力于游戏 AI, 并撰写了大量有关该主题的论文和出版物。

他现在通过承包的方式从事游戏 AI 工作，同时是多本游戏 AI 图书的撰稿人和技术编辑。

## Thomas Young

---

thomas.young@bigfoot.com

Thomas Young 是一名自由程序员，在位于法国南部阳光明媚的家中办公。他当前为比利时的 Appeal 公司工作，负责为 *Outcast 2* 开发寻径系统。成为自由程序员之前，他在谢菲尔德的 Gremlin Interactive 工作过多年，在多个项目中承担 AI 寻径方面的编码工作，其中最著名的项目是 *Soulbringer*。他的少年时光是在为 Amiga 芯片组编写硬件代码中度过的，随后他前往斯塞克斯大学攻读人工智能学士学位。

## Michael Zarozinski

---

michaelz@louderthanabomb.com

www.louderthanabomb.com

Michael Zarozinski 是 Louder Than A Bomb! Software——一家致力于将 AI 用于娱乐软件的公司——的老板和负责人。Michael 终生迷恋人工智能。他编写过 TI-99/4A 版本的 *Eliza*。有关 Michael 的详细信息，请访问 LouderThanABomb! 网站，其中的内容可能比您想知道的还要多。Michael 在业余时间管理网站 aiGuru.com，其中包括各种有关 AI 的知识以及他对 AI 和非 AI 主题的见解。



# 目 录

## 第 1 章 通用编程技术

绪 论 .....	3
<i>Scott Bilas</i>	
1.1 优化 C++ 游戏 .....	5
<i>Andrew Kirmse</i>	
1.1.1 对象的创建和销毁 .....	5
1.1.2 内存管理 .....	8
1.1.3 虚拟函数 .....	9
1.1.4 代码长度 .....	10
1.1.5 标准模板库 .....	11
1.1.6 高级特性 .....	12
1.1.7 参考文献 .....	13
1.2 内联函数和宏 .....	14
<i>Peter Dalton</i>	
1.2.1 内联函数的优点 .....	14
1.2.2 何时使用内联函数 .....	15
1.2.3 何时使用宏 .....	16
1.2.4 微软特有的情况 .....	16
1.2.5 参考文献 .....	17
1.3 抽象接口编程 .....	18
<i>Noel Liopis</i>	
1.3.1 抽象接口 .....	18
1.3.2 添加一个工厂 (factory) .....	19
1.3.3 抽象接口特性 .....	21
1.3.4 一切都是代价的 .....	23
1.3.5 结论 .....	24
1.3.6 参考文献 .....	24
1.4 从 DLL 中导出 C++ 类 .....	25
<i>Herb Marselas</i>	
1.4.1 导出函数 .....	25
1.4.2 导出类 .....	25
1.4.3 导出类成员函数 .....	27
1.4.4 导出虚拟类成员函数 .....	27

1.4.5 总结 .....	28
<b>1.5 避免 DLL 困境</b> .....	29
<i>Herb Marselas</i>	
1.5.1 显式链接还是隐式链接 .....	29
1.5.2 LoadLibrary 和 GetProcAddress .....	30
1.5.3 提防 DirectX .....	30
1.5.4 使用操作系统特有的特性 .....	31
1.5.5 总结 .....	32
<b>1.6 动态类型信息</b> .....	34
<i>Scott Wakeling</i>	
1.6.1 动态类型信息类简介 .....	34
1.6.2 暴露和查询 DTI .....	35
1.6.3 继承的含义是“是一个” .....	36
1.6.4 处理通用对象 .....	37
1.6.5 实现永久性类型信息 .....	38
1.6.6 将永久性类型信息用于游戏保存数据库中 .....	39
1.6.7 结论 .....	40
1.6.8 参考文献 .....	40
<b>1.7 用于通用 C++ 成员访问的属性类</b> .....	41
<i>Charles Cafrelli</i>	
1.7.1 代码 .....	41
1.7.2 其他用途 .....	44
1.7.3 推荐读物 .....	44
<b>1.8 一个游戏实体工厂</b> .....	45
<i>François Dominic Laramée</i>	
1.8.1 组件 .....	45
1.8.2 flyweight 类、行为类和导出类 .....	45
1.8.3 flyweight 对象 .....	46
1.8.4 SAMMy, 你在哪里? .....	46
1.8.5 行为类层次 .....	47
1.8.6 使用模板方法模式来完成行为任务 .....	48
1.8.7 导出类 .....	49
1.8.8 实体工厂 .....	50
1.8.9 在运行阶段选择策略 .....	52
1.8.10 最后的注意事项 .....	53
1.8.11 参考文献 .....	53
<b>1.9 在 C++ 添加摒弃功能</b> .....	55
<i>Noel Llopis</i>	
1.9.1 可能的解决方案 .....	55



1.9.2	理想的解决方案 .....	55
1.9.3	使用和指定被摒弃的函数 .....	56
1.9.4	使用 C++ 实现摒弃功能 .....	56
1.9.5	可改进的地方 .....	58
1.9.6	致谢 .....	58
1.9.7	参考文献 .....	58
<b>1.10</b>	<b>一个插入式调试内存管理器 .....</b>	<b>59</b>
	<i>Peter Dalton</i> .....	
1.10.1	内存管理器初步 .....	59
1.10.2	内存管理器的记录工作 .....	60
1.10.3	报告信息 .....	62
1.10.4	注意事项 .....	63
1.10.5	进一步的改进 .....	64
1.10.6	参考文献 .....	64
<b>1.11</b>	<b>一个内置的游戏剖析模块 .....</b>	<b>66</b>
	<i>Jeff Evertt</i> .....	
1.11.1	有关剖析的基本知识 .....	66
1.11.2	商用工具 .....	67
1.11.3	为何要自己开发模块 .....	67
1.11.4	剖析模块 (Profiling module) 的需求 .....	68
1.11.5	架构和实现 .....	68
1.11.6	实现的细节 .....	69
1.11.7	分析数据 .....	69
1.11.8	有关实现的注意事项 .....	70
<b>1.12</b>	<b>用于 Windows 游戏的线性编程模型 .....</b>	<b>71</b>
	<i>Javier F. Otaegui</i> .....	
1.12.1	更新背景 .....	71
1.12.2	解决方案: 多线程 (Multithreading) .....	72
1.12.3	参考文献 .....	75
<b>1.13</b>	<b>栈缠绕 .....</b>	<b>76</b>
	<i>Bryon Hapgood</i> .....	
1.13.1	简单的 TempRet .....	76
1.13.2	TempRet 链 .....	77
1.13.3	Thunking .....	78
1.13.4	递归 .....	80
<b>1.14</b>	<b>自我修改的代码 .....</b>	<b>82</b>
	<i>Bryon Hapgood</i> .....	
1.14.1	RAM 代码的原理 .....	82
1.14.2	一个快速的 Bit Blitter .....	83

<b>1.15 使用资源文件来管理文件</b>	91
<i>Bruno Sousa</i>	
1.15.1 何为资源文件	91
1.15.2 设计	92
1.15.3 实现	93
1.15.4 有关实现的最后一些说明	94
1.15.5 结论	95
1.15.6 参考文献	95
<b>1.16 游戏输入的记录和重放</b>	96
<i>Bruce Dawson</i>	
1.16.1 记录输入有何用途	96
1.16.2 原理	97
1.16.3 测试输入记录功能	100
1.16.4 结论	100
1.16.5 参考文献	101
<b>1.17 一个灵活的文本分析系统</b>	102
<i>James Boer</i>	
1.17.1 分析系统	102
1.17.2 宏、头文件和预处理技术	103
1.17.3 该分析系统的结构	104
1.17.4 小结	106
<b>1.18 一个通用的调节器</b>	107
<i>Lasse Staff Jensen</i>	
1.18.1 需求分析	107
1.18.2 实现	107
1.18.3 使用	112
1.18.4 图形用户界面	112
1.18.5 附注	114
1.18.6 致谢	114
<b>1.19 生成真正的随机数</b>	115
<i>Pete Isensee</i>	
1.19.1 伪随机	115
1.19.2 真正随机	115
1.19.3 随机输入源	116
1.19.4 硬件源	116
1.19.5 混合函数	117
1.19.6 局限性	117
1.19.7 实现	117
1.19.8 GenRand 的随机程度	119

1.19.9 参考文献 .....	119
<b>1.20 使用 Bloom 过滤器来提高计算性能 .....</b>	<b>120</b>
<i>Mark Fischer</i>	
1.20.1 Bloom 的方式 .....	120
1.20.2 可能的情形 .....	120
1.20.3 工作原理 .....	121
1.20.4 定义 .....	121
1.20.5 范例 1 .....	121
1.20.6 范例 2 .....	125
1.20.7 最后的说明 .....	125
1.20.8 结论 .....	126
1.20.9 参考文献 .....	126
<b>1.21 3Ds MAX 中的 Skin 导出器和动画工具包 .....</b>	<b>127</b>
<i>Marco Tombesi</i>	
1.21.1 导出 .....	128
1.21.2 参考文献 .....	137
<b>1.22 在视频游戏中使用 Web 摄像机 .....</b>	<b>138</b>
<i>Nathan d'Obrenan</i>	
1.22.1 初始化 Web 摄像机捕获窗口 .....	138
1.22.2 操纵 Web 摄像机数据 .....	143
1.22.3 结论 .....	147
1.22.4 参考文献 .....	147

## 第 2 章 数学技巧

绪 论 .....	151
<i>Eddie Edwards</i>	
<b>2.1 浮点计算技巧：使用 IEEE 浮点格式以提高性能 .....</b>	<b>152</b>
<i>Yossarian King</i>	
2.1.1 概述 .....	152
2.1.2 IEEE 浮点格式 .....	152
2.1.3 浮点数技巧 .....	153
2.1.4 用于正弦和余弦函数的线性查找表 .....	158
2.1.5 平方根函数的对数优化 .....	160
2.1.6 优化任何函数 .....	161
2.1.7 性能测量 .....	163
2.1.8 结论 .....	163
2.1.9 参考文献 .....	164
<b>2.2 矢量和平面技巧 .....</b>	<b>165</b>
<i>John Olsen</i>	

2.2.1	相对于碰撞面的高度 .....	165
2.2.2	找出碰撞点 .....	166
2.2.3	到碰撞点的距离 .....	167
2.2.4	反射式碰撞 .....	168
2.2.5	阻尼碰撞 .....	170
<b>2.3</b>	<b>一种快速、健壮的计算 3D 线段交点的方法 .....</b>	<b>172</b>
	<i>Graham Rhodes</i>	
2.3.1	这种算法健壮的原因 .....	172
2.3.2	问题描述 .....	172
2.3.3	推导闭式解决方案方程 .....	174
2.3.4	线段 .....	179
2.3.5	实现描述 .....	181
2.3.6	可优化的地方 .....	181
2.3.7	结论 .....	182
2.3.8	参考文献 .....	182
<b>2.4</b>	<b>反向弹道计算 .....</b>	<b>183</b>
	<i>Aaron Nicholls</i>	
2.4.1	一种特殊情况 .....	184
2.4.2	优化实现 .....	189
2.4.3	总结 .....	190
<b>2.5</b>	<b>平行移动镜头 .....</b>	<b>191</b>
	<i>Carl Dougan</i>	
2.5.1	技术 .....	191
2.5.2	结论 .....	194
2.5.3	参考文献 .....	194
<b>2.6</b>	<b>平滑的基于四元数的 <math>C^2</math> 飞行路径 .....</b>	<b>195</b>
	<i>Alex Vlachos</i>	
2.6.1	导论 .....	195
2.6.2	位置插值 .....	195
2.6.3	朝向插值 .....	196
2.6.4	旋转方向和选择性求负 .....	197
2.6.5	四元数样条线插值 .....	198
2.6.6	有理映射中的奇异点 .....	199
2.6.7	镜头剪接 .....	199
2.6.8	代码 .....	199
2.6.9	参考文献 .....	200
<b>2.7</b>	<b>递归逐维分组：一种快速的碰撞检测算法 .....</b>	<b>201</b>
	<i>Steve Rabin</i>	
2.7.1	其他应用 .....	202

2.7.2	该算法的缺陷 .....	205
2.7.3	查找碰撞物体对 .....	206
2.7.4	时间复杂度 .....	208
2.7.5	结论 .....	209
2.7.6	参考文献 .....	209
<b>2.8</b>	<b>不规则碎片编程 .....</b>	<b>210</b>
	<i>Jesse Laeuchli</i>	
2.8.1	无规则碎片 .....	210
2.8.2	断层无规则碎片 .....	211
2.8.3	FBM .....	211
2.8.4	实现 .....	212
2.8.5	使用 FBM .....	215
2.8.6	参考文献 .....	216

### 第3章 人工智能

绪 论.....	219
<i>Steve Rabin</i>	
<b>3.1 AI 优化策略.....</b>	<b>221</b>
<i>Steve Rabin</i>	
3.1.1 策略 1: 使用事件驱动行为而非轮询 .....	221
3.1.2 策略 2: 减少重复计算 .....	222
3.1.3 策略 3: 由管理员集中进行协调 .....	222
3.1.4 策略 4: 不那么频繁地运行 AI.....	222
3.1.5 策略 5: 将处理工作分散到多帧中完成 .....	223
3.1.6 策略 6: 利用细节级 AI.....	223
3.1.7 策略 7: 只解决问题的一部分 .....	223
3.1.8 策略 8: 离线完成困难的工作 .....	224
3.1.9 策略 9: 使用突发行为以避免编写脚本 .....	224
3.1.10 策略 10: 通过连续记录来分摊查询成本 .....	224
3.1.11 策略 11: 重新考虑问题.....	225
3.1.12 结 论 .....	225
3.1.13 参考文献 .....	226
<b>3.2 用于游戏对象 AI 的微线程 .....</b>	<b>227</b>
<i>Bruce Dawson</i>	
3.2.1 一个更简单的方法 .....	228
3.2.2 微线程 .....	229
3.2.3 栈管理 .....	230
3.2.4 并发症 .....	231



3.2.5	结论 .....	232
3.2.6	参考文献 .....	232
<b>3.3</b>	<b>使用微线程管理 AI .....</b>	<b>233</b>
	<i>Simon Carter</i>	
3.3.1	拼凑 .....	233
3.3.2	良好的行为 .....	233
3.3.3	了然于胸 .....	234
3.3.4	并发症 .....	236
3.3.5	结论 .....	238
3.3.6	参考文献 .....	239
<b>3.4</b>	<b>一种 RTS 命令排队体系结构 .....</b>	<b>240</b>
	<i>Steve Rabin</i>	
3.4.1	RTS 命令 .....	240
3.4.2	命令排队 .....	240
3.4.3	循环命令 .....	242
3.4.4	结论 .....	244
3.4.5	参考文献 .....	244
<b>3.5</b>	<b>一种基于分片的高性能视域和搜索系统 .....</b>	<b>245</b>
	<i>Matt Pritchard</i>	
3.5.1	概述 .....	245
3.5.2	定义 .....	245
3.5.3	组件 1: 各个玩家的可见性地图 .....	246
3.5.4	组件 2: LOS 模板 .....	246
3.5.5	组件 3: 合并的可视性地图 .....	248
3.5.6	改进搜索 .....	249
3.5.7	结论 .....	251
<b>3.6</b>	<b>创建影响力地图 .....</b>	<b>252</b>
	<i>Paul Tozour</i>	
3.6.1	影响力地图 .....	252
3.6.2	一个简单的影响力地图 .....	253
3.6.3	影响力地图中的单元格数据 .....	254
3.6.4	计算合意值 .....	255
3.6.5	确定最佳的单元格大小 .....	256
3.6.6	影响力传播 .....	256
3.6.7	考虑地形 .....	257
3.6.8	特别考虑 .....	258
3.6.9	刷新影响力地图 .....	259
3.6.10	3D 环境中的影响力地图 .....	259
3.6.11	参考文献和推荐读物 .....	260

**3.7 策略评估技术** ..... 261

*Paul Tozour*

3.7.1 资源分配树 ..... 261

3.7.2 计算希望的资源分配 ..... 262

3.7.3 判断当前的分配情况 ..... 263

3.7.4 策略决策 ..... 263

3.7.5 值的估量 ..... 264

3.7.6 依存图 ..... 264

3.7.7 依存图节点 ..... 265

3.7.8 经济规划 ..... 265

3.7.9 查找脆弱的依存 ..... 266

3.7.10 策略推理 ..... 266

3.7.11 玩家个性..... 267

3.7.12 总结 ..... 267

3.7.13 参考文献 ..... 267

**3.8 3D 游戏中的地形推理** ..... 269

*William van der Sterren*

3.8.1 以方便推理的方式表示地形 ..... 269

3.8.2 中继点 (waypoint) ..... 270

3.8.3 范例地形和 AI 需求 ..... 270

3.8.4 战术分析 (tactical analysis) ..... 270

3.8.5 从战术价值到中继点属性 ..... 271

3.8.6 计算中继点属性 ..... 272

3.8.7 从经验中学习 ..... 274

3.8.8 将地形推理加入到游戏中 ..... 275

3.8.9 其他应用 ..... 275

3.8.10 结论 ..... 276

3.8.11 参考文献和推荐读物 ..... 276

**3.9 用于可视点寻径的扩展几何体** ..... 277

*Thomas Young*

3.9.1 定义碰撞模型 ..... 277

3.9.2 多边形寻径 ..... 278

3.9.3 扩展并解决问题 ..... 278

3.9.4 凸多边形的闵可夫斯基和 ..... 279

3.9.5 扩展非凸几何体 ..... 280

3.9.6 选择碰撞形状 ..... 281

3.9.7 结论 ..... 281

3.9.8 参考文献 ..... 281

**3.10 优化可视点寻径** ..... 283

*Thomas Young*

3.10.1 可视点寻径 (points-of-visibility pathfinding) ..... 283

3.10.2 存储到每个点的最短路径 ..... 283

3.10.3 将凸角相连 ..... 284

3.10.4 轮廓区 ..... 285

3.10.5 将轮廓区用于空间分区系统 ..... 287

3.10.6 结论 ..... 287

3.10.7 参考文献 ..... 287

**3.11 有齿物群的模拟：捕食者和猎物** ..... 288

*Steven Woodcock*

3.11.1 全新的世界 ..... 289

3.11.2 有齿物群的模拟 ..... 291

3.11.3 局限性和可改进的地方 ..... 291

3.11.4 参考文献 ..... 292

**3.12 一个用 C++编写的通用模糊状态机** ..... 293

*Eric Dybsand*

3.12.1 为何在游戏中使用 FuSM ..... 294

3.12.2 如何在游戏中使用 FuSM ..... 294

3.12.3 复习《游戏编程精粹 1》中的 C++通用有限状态机 ..... 295

3.12.4 将通用 FSM 修改为 FuSM ..... 295

3.12.5 在游戏中使用模糊逻辑 ..... 296

3.12.6 参考文献 ..... 296

**3.13 避免模糊系统中的组合激增** ..... 297

*Michael Zaroinski*

3.13.1 问题 ..... 297

3.13.2 解决方案 ..... 298

3.13.3 范例 ..... 299

3.13.4 结论 ..... 303

3.13.5 参考文献 ..... 303

**3.14 一个在游戏中使用神经网络的例子** ..... 304

*John Manslow*

3.14.1 游戏 ..... 304

3.14.2 多玩家感知器 ..... 304

3.14.3 选择输入 ..... 306

3.14.4 收集数据 ..... 306

3.14.5 训练 MLP ..... 307

3.14.6 结果 ..... 308

3.14.7 结论 ..... 308

3.14.8 参考文献 .....	309
-------------------	-----

## 第4章 几何体管理

绪 论.....	313
<i>Eric Lengyel</i>	
<b>4.1 各种 VIPM 方法的比较</b> .....	314
<i>Tom Forsyth</i>	
4.1.1 考虑因素 .....	314
4.1.2 Vanilla VIPM .....	316
4.1.3 跳带 .....	319
4.1.4 多层跳带 .....	319
4.1.5 混合模式 VIPM .....	320
4.1.6 混合模式跳带 .....	321
4.1.7 滑窗 .....	321
4.1.8 小结 .....	324
4.1.9 参考文献 .....	325
<b>4.2 使用联锁分片简化地形</b> .....	326
<i>Greg Snook</i>	
4.2.1 分片的重访问 .....	327
4.2.2 生成地图 .....	328
4.2.3 分片模板 .....	328
4.2.4 消除难看的接缝 .....	329
4.2.5 更好、更快、更强 .....	330
4.2.6 结论 .....	331
4.2.7 参考文献 .....	331
<b>4.3 快速可视剔除、射线跟踪以及范围搜索的球形树</b> .....	332
<i>John W. Ratcliff</i>	
4.3.1 包围球 .....	332
4.3.2 使用球形树 .....	332
4.3.3 演示应用程序 .....	333
<b>4.4 压缩的轴向包围盒树</b> .....	335
<i>Miguel Gomez</i>	
4.4.1 概览层次排序方法 .....	335
4.4.2 AABB 树 .....	336
4.4.3 构建 AABB 树 .....	336
4.4.4 压缩 AABB 树 .....	337
4.4.5 近似范围 .....	337
4.4.6 利用冗余 .....	338
4.4.7 运行时效 .....	339

4.4.8 将来的工作 .....	339
4.4.9 参考文献 .....	339
<b>4.5 直接访问四叉树查找 .....</b>	<b>340</b>
<i>Matt Pritchard</i>	
4.5.1 性能剖析 .....	340
4.5.2 消除中间阻碍 .....	341
4.5.3 条件和要求 .....	341
4.5.4 判断树层 .....	341
4.5.5 位置映射 .....	343
4.5.6 判断位置 .....	344
4.5.7 遍历四叉树 .....	344
4.5.8 优化四叉树 .....	344
<b>4.6 近似鱼缸折射 .....</b>	<b>347</b>
<i>Alex Vlachos</i>	
4.6.1 鱼缸观察现象 .....	347
4.6.2 提高其真实性 .....	349
4.6.3 结论 .....	349
<b>4.7 渲染打印分辨率的屏幕快照 .....</b>	<b>350</b>
<i>Alex Vlachos</i>	
4.7.1 基本算法 .....	350
4.7.2 忠告及注意 .....	353
4.7.3 结论 .....	353
4.7.4 参考文献 .....	353
<b>4.8 对任意表面应用贴花 .....</b>	<b>354</b>
<i>Eric Lengyel</i>	
4.8.1 算法 .....	354
4.8.2 三角形剪裁 .....	356
4.8.3 实现代码 .....	357
4.8.4 参考文献 .....	357
<b>4.9 用天空包围盒渲染远景 .....</b>	<b>358</b>
<i>Jason Shankel</i>	
4.9.1 基本技术 .....	358
4.9.2 天空包围盒分辨率 .....	359
4.9.3 天空包围盒大小 .....	359
4.9.4 渲染场景 .....	360
4.9.5 立方体环境映射 .....	360
4.9.6 生成天空包围盒纹理 .....	361
4.9.7 结论 .....	361
4.9.8 源代码 .....	361

<b>4.10 自阴影角色</b> .....	362
<i>Alex Vlachos, David Gosselin, Jason L. Mitchell</i>	
4.10.1 研究回顾 .....	362
4.10.2 角色几何分割 .....	362
4.10.3 渲染纹理 .....	362
4.10.4 渲染角色 .....	363
4.10.5 结论 .....	364
4.10.6 参考文献 .....	364
<b>4.11 经典的 <i>Super Mario 64</i> 游戏第三人称控制和动画</b> .....	365
<i>Steve Rabin</i>	
4.11.1 设置 .....	365
4.11.2 转换控制器的输入 .....	365
4.11.3 旋转角色 .....	367
4.11.4 角色移动 .....	368
4.11.5 角色动画 .....	368
4.11.6 <i>Super Mario 64</i> 动画分析 .....	370
4.11.7 结论 .....	371
4.11.8 参考文献 .....	371

## 第 5 章 图形显示

绪 论 .....	375
<i>D. Sim Dietrich Jr</i>	
<b>5.1 卡通渲染：实时轮廓边缘检测与渲染</b> .....	376
<i>Carl S. Marshall</i>	
5.1.1 着墨器 (Inker) .....	376
5.1.2 重要的边 .....	377
5.1.3 轮廓边缘检测技术 .....	377
5.1.4 基于边的着墨 .....	377
5.1.5 可编程顶点着色器着墨 .....	379
5.1.6 高级纹理特征着墨 .....	381
5.1.7 结论 .....	381
5.1.8 参考文献 .....	382
<b>5.2 使用纹理映射的卡通渲染与可编程顶点着色器</b> .....	383
<i>Adam Lake</i>	
5.2.1 卡通着色技术 .....	383
5.2.2 上色 .....	383
5.2.3 可编程顶点着色器 .....	386
5.2.4 结论 .....	388
5.2.5 参考文献 .....	388



<b>5.3 动态逐像素光照技术</b> .....	389
<i>Dan Ginsburg, Dave Gosselin</i>	
5.3.1 动态光照贴图的 3D 纹理 .....	389
5.3.2 Dot3 凹凸贴图 (Bump Mapping) .....	391
5.3.3 使用立方贴图规一化 .....	395
5.3.4 逐像素聚光灯 (Per-Pixel Spotlight) .....	396
5.3.5 参考文献 .....	396
<b>5.4 使用 3D 硬件生成过程云彩</b> .....	398
<i>Kim Pallister</i>	
5.4.1 云彩性质 .....	398
5.4.2 生成随机数 .....	399
5.4.3 噪音多个倍频的动画 .....	401
5.4.4 贴图到天空几何体 .....	403
5.4.5 功能延伸 .....	404
5.4.6 硬件限制 .....	404
5.4.7 可伸缩性 .....	405
5.4.8 结论 .....	405
5.4.9 参考文献 .....	406
<b>5.5 针对较快镜头眩光的纹理屏蔽</b> .....	407
<i>Chris Maughan</i>	
5.5.1 镜头眩光遮挡 .....	407
5.5.2 硬件问题 .....	407
5.5.3 纹理屏蔽 .....	409
5.5.4 性能考虑 .....	410
5.5.5 改进 .....	411
5.5.6 示例代码 .....	411
5.5.7 替代途径 .....	412
5.5.8 参考文献 .....	412
<b>5.6 实用优先缓冲阴影</b> .....	413
<i>D. Sim Dietrich Jr.</i>	
5.6.1 比较优先缓冲与深度缓冲 .....	415
5.6.2 解决锯齿化问题 .....	416
5.6.3 混合途径 .....	417
5.6.4 小结 .....	418
5.6.5 参考文献 .....	418
<b>5.7 替用体技术：添加点缀</b> .....	419
<i>Tom Forsyth</i>	
5.7.1 整个过程 .....	419
5.7.2 渲染替换体 .....	420

5.7.3 更新试探法 ..... 423

5.7.4 效率 ..... 424

5.7.5 预测 ..... 424

5.7.6 小结 ..... 425

**5.8 硬件加速过程纹理动画中的运算** ..... 426

*Greg James*

5.8.1 硬件运算 ..... 426

5.8.2 将来的工作 ..... 435

5.8.3 致谢 ..... 435

5.8.4 示例源码 ..... 435

5.8.5 参考文献 ..... 435

第 6 章 音频编程

绪 论 ..... 439

*James Boer*

**6.1 游戏音频设计模式** ..... 440

*Scott Patterson*

6.1.1 桥接 (Bridge) ..... 440

6.1.2 外观 (Façade) ..... 441

6.1.3 合成 (Composite) ..... 441

6.1.4 代理 (Proxy) ..... 442

6.1.5 修饰器 (Decorator) ..... 442

6.1.6 命令 (Command) ..... 443

6.1.7 备忘录 (Memento) ..... 443

6.1.8 观测器 (Observer) ..... 443

6.1.9 大泥球 (Big Ball of Mud) (也称做“意大利面条式”代码) ..... 444

6.1.10 结论 ..... 445

6.1.11 参考文献 ..... 445

**6.2 在采样合成器中声音的同步重用技术** ..... 446

*Thomas Engel*

6.2.1 存在的问题 ..... 446

6.2.2 解决方案的思路 ..... 447

6.2.3 解决方案 ..... 447

6.2.4 结论 ..... 448

**6.3 软件 DSP 效果** ..... 450

*Ian Lewis*

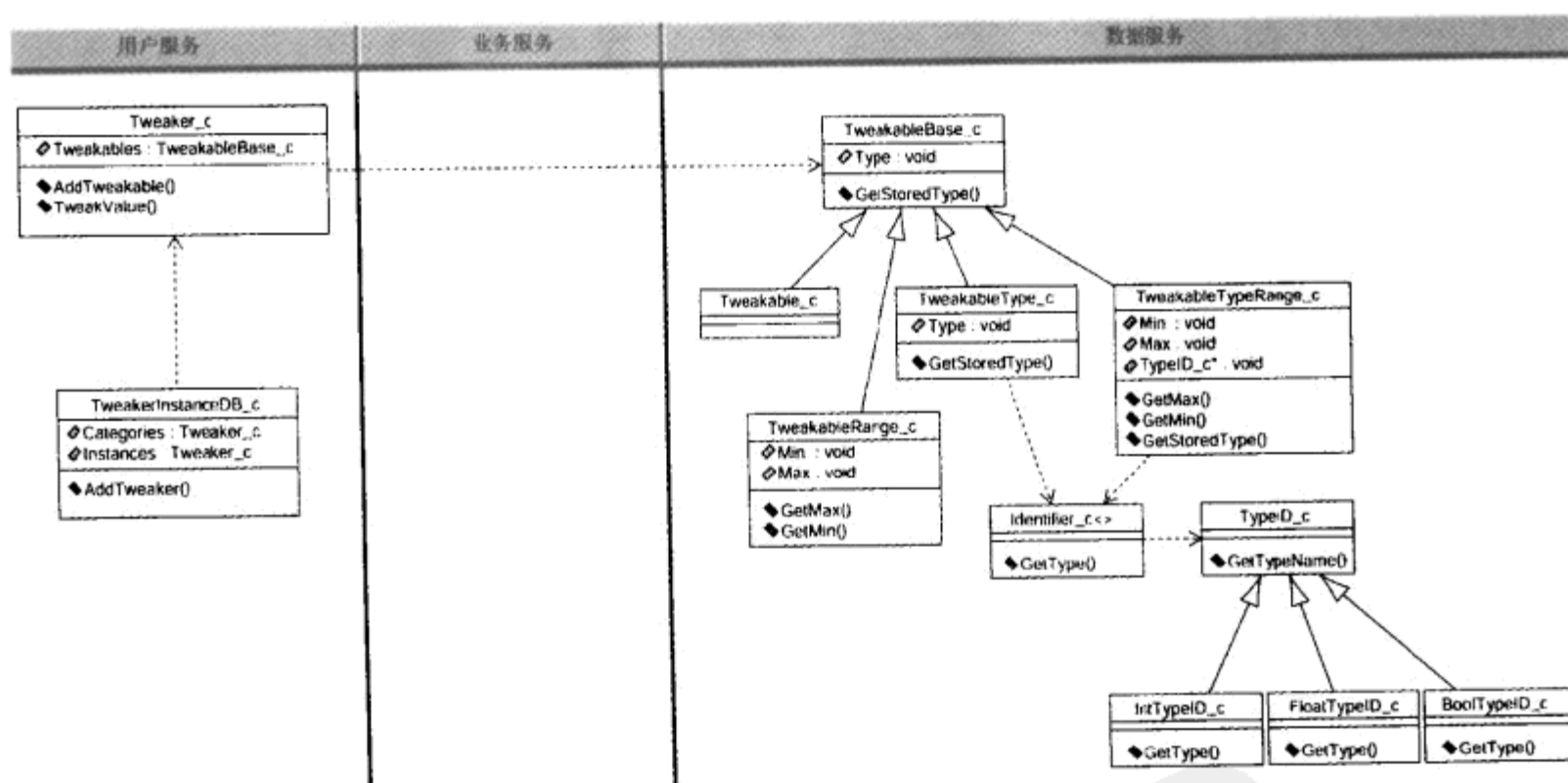
6.3.1 滤波 ..... 450

6.3.2 卷积 (convolution) ..... 451

6.3.3 延迟 ..... 451

6.3.4	插值 (interpolation) .....	452
6.3.5	参考文献 .....	453
<b>6.4</b>	<b>数字音频的交互式处理管线 .....</b>	<b>454</b>
	<i>Keith Weiner</i>	
6.4.1	简介 .....	454
6.4.2	讨论 .....	456
6.4.3	代码 .....	458
6.4.4	额外注释 .....	461
6.4.5	结论 .....	462
<b>6.5</b>	<b>游戏中的基本音乐音序器 .....</b>	<b>463</b>
	<i>Scott Patterson</i>	
6.5.1	音乐流与音序 .....	463
6.5.2	核心计算机音乐概念 .....	464
6.5.3	计算机音序器实现 .....	466
6.5.4	音频合成 (composite) 控制 .....	472
6.5.5	源码 .....	473
6.5.6	结论 .....	473
6.5.7	参考文献 .....	473
<b>6.6</b>	<b>用于游戏的交互式音序器 .....</b>	<b>474</b>
	<i>Scott Patterson</i>	
6.6.1	音乐联想 .....	474
6.6.2	音乐意义 .....	474
6.6.3	过渡 .....	476
6.6.4	过渡类型 .....	476
6.6.5	控制粒度 .....	477
6.6.6	目标控制 .....	477
6.6.7	设计示例 .....	479
6.6.8	源码 .....	480
6.6.9	结论 .....	480
6.6.10	参考文献 .....	480
<b>6.7</b>	<b>底层声音 API .....</b>	<b>481</b>
	<i>Ian Lewis</i>	
	核心类 .....	481
	索引 .....	483

## 通用编程技术





## 绪 论

Scott Bilas, Gas Powered Games

scottb@aa.net

**阅**读本章内容之前，请读者回答这样一个问题：何为工程师（或者说您从事的是什么样的工作）？工程师是编写代码或能够使得大量三角形快速地绘制出来的人吗？当然不是。很多工程师并不编写大量代码。工程学研究的是如何解决问题。

那么，在游戏开发工作室，工程师扮演的是什么角色呢？答案很简单：工程师为开发小组服务。撇开人们常常同时扮演多种角色不谈，开发小组通常分为三个主要的小组。首先是设计人员，他们负责设计游戏的脉络，并确保其有趣；其次是内容开发人员——通常是艺术家、作曲家和关卡设计师，他们负责制作设计方案所需的素材；最后是工程师，他们负责实现设计方案——所需的时间通常要长得多。工程师负责创建一个系统，让玩家能够同游戏中的素材交互。换句话说，工程师通过创建工具，来为小组提供服务。服务？工具？

人们通常认为是工具的（如关卡编辑器或脚本批处理器）只是画面中可见的、突出的部分。工程师创建的所有东西都是工具，无论是用于播放或混合艺术家创建的动画的引擎、用于管理交互式内容的游戏内嵌的数据库，还是跟踪内存使用情况的调试设施。我们的终极目标是，创建出能够实现游戏设计方案的软件工具。

我们创建的很多工具并没有嵌入到游戏中，它们用于创建和处理数据，以便高效地播放或模拟游戏引擎。这样的工具包括 BSP 编译器、纹理压缩程序、资源组合器和关卡编辑器。这些工具通常只供开发人员和模型设计人员（mod maker）使用。没有它们，游戏也能运行；但对于创建游戏而言，它们是必不可少的。当前，这种工具开始被集成到游戏引擎中。这样做旨在加快游戏的开发进度，因为这使得游戏中很大一部分内容可以在游戏运行的同时进行编辑。

我们创建的其他工具用于游戏中，其中常见的有脚本引擎、物理线现象模拟器和用户界面。所有这些工具的核心都是一个命令处理器，它接收指令并根据这些信息完成一些有用的工作。换句话说，这些工具按设计方案的要求对内容进行解释和塑造，从而实现可供玩家体验的游戏。

有时候，工程师为帮助完成自己的工作，也需要创建工具，这些工具被嵌入到游戏中。现代的游戏非常复杂，为高效地调试和优化它们，游戏本身必须提供大量的支持功能。由于机器中安装的内存有几十兆（不久后



将为几百兆)，因此游戏中包含的内容多得令人难以置信，为免遭灭顶之灾，我们必须创建复杂的支持系统。我们创建调试系统、剖析程序（`profiler`）、类型数据库、实用程序库、栈转储（`stack-dumping`）异常处理程序、记录工具、自检代码等，真是数不胜数。本章所介绍的正是这样的工具，这些文章将帮助读者打下坚实的游戏开发基础。祝您阅读愉快！



## 1.1 优化 C++ 游戏

Andrew Kirmse, LucasArts Entertainment Company  
ark@alum.mit.edu

与用 C 语言编写的游戏相比，编写良好的 C++ 游戏通常更容易维护，可重用性也更好，但使用 C++ 来编写游戏值得吗？复杂的 C++ 程序在运行速度上能够同 C 程序媲美吗？

只要使用优秀的编译器，并对 C++ 非常熟悉，使用这种语言完全可以创建出高效的游戏。本文介绍可用于提高游戏速度的技术。这里假设读者已确信使用 C++ 的好处，并熟悉常规的优化原理（可参见参考文献）。

这里有必要重申的一种通用规则是，剖析（profile）是至关重要的。在剖析时，程序员常犯两种错误。首先，没有对合适的代码进行优化。程序的大部分代码对性能的影响并不大，对这些代码进行优化无疑是浪费时间。有关哪些代码对性能的影响至关重要的直觉是不可靠的，只有通过直接测试才能确定这一点。其次，有时候程序员对代码进行优化时，实际上会使代码的运行速度更慢。对于 C++ 代码尤其如此，在 C++ 中，一行看起来非常简单的代码可能生成大量的机器代码。程序员应查看编译器的输出，并经常对其进行剖析。

### 1.1.1 对象的创建和销毁

对象的创建和销毁是 C++ 中的一个核心概念，编译器常常因此为您生成大量代码。设计不佳的程序将花大量的时间来调用构造函数、复制对象以及生成开销高昂的临时对象。幸运的是，一些常识和简单的规则，能够使得使用大量对象的代码运行得和 C 代码几乎一样快。

- 将对象的创建工作推迟到需要时进行

不运行代码时，其速度是最快的；因此，对于无须使用的对象，不应创建它。请看下面的代码中：

```
void Function(int arg)
{
    Object obj;
    if (arg == 0)
        return;
    ...
}
```

这段代码中，即使 `arg` 为零，也将调用对象的构造函数和析构函数。如果 `arg` 总是为零，且对象本身将分配内存，这种浪费将剧增。解决之道是将 `obj` 声明语句移到 `if` 语句的后面。

然而，在循环中声明并非不重要的对象时应小心。如果在循环中，将对象的创建工作推迟到需要时进行，则每次迭代都将创建和销毁该对象。因此，最好在循环之前声明对象，这样将只创建和销毁该对象一次。如果需要在内嵌的循环中调用一个函数，而该函数在堆栈中创建一个对象，则应在循环的外面创建该对象，并通过引用将其传递给函数。

- 使用初始化列表

请看下面的类：

```
class Vehicle
{
public:
    Vehicle(const std::string &name) // Don't do this!
    {
        mName = name;
    }
private:
    std::string mName;
};
```

由于成员变量将在构造函数体调用之前创建，因此上述代码将为字符串 `mName` 调用构造函数，然后调用操作符`=`复制创建的字符串。对于上述范例，尤其糟糕的是，缺省的字符串构造函数可能分配大量的内存——比用于存储 `Vehicle` 的构造函数赋给该变量的名称所需的内存多。下面的代码将好得多，且不会调用操作符`=`。另外，如果知道更详细的信息（这里为实际要存储的字符串），则非缺省的字符串构造函数的效率通常更高，同时如果 `Vehicle` 的构造函数体为空，则编译器将对代码进行优化，而不会调用该函数。

```
class Vehicle
{
public:
    Vehicle(const std::string &name) : mName(name)
    {
    }
private:
    std::string mName;
};
```

- 先递增而不是后递增

代码 `x = y++` 存在的问题在于，递增函数必须先复制 `y` 的初始值，将 `y` 加 1，然后再返回 `y` 的初始值。因此后递增需要创建一个临时对象，而先递增不需要。就整型而言，这不会带来额外的开销；但对于用户定义的类型，这将是经济的。因此，在可能的情况下，应采用先递增的方式；在 `for` 循环的迭代器中，通常可以这样做。

- 避免操作符按值返回结果

在 C++ 中，矢量相加的规范代码如下：

```
Vector operator+(const Vector &v1, const Vector &v2)
```

该操作符将返回一个新的 `Vector` 对象，且按值返回该对象。虽然这样使得我们可以编写诸如 `v = v1 + v2` 等可读性较强的表达式，但如果需要大量进行矢量相加操作，则创建临时 `Vector` 对象并复制它们的开销将太大。有时候，可以对代码进行改编，让编译器能够对代码进行优化，从而避免创建临时对象（这被称为“返回值优化”）。不过，通常编写稍微难看些，但速度更快的下述代码更佳：

```
void Vector::Add(const Vector &v1, const Vector &v2)
```

请注意，操作符`+=`不存在这样的问题，因为它修改第一个参数，而不用返回一个临时对象。因此，在可能的情况下，应使用操作符`+=`，而不是`+`。

- 使用轻量级构造函数

在上述范例中，`Vector` 类的构造函数应将其元素初始化为零吗？在有些情况下，这样做可能有用，但这将使得每个调用者都需要对其进行初始化，而不管他是否将使用这些元素。另外，临时矢量及其成员变量将隐式地带来其他额外的开销。

优秀的编译器将把一些多余的代码优化掉，但为何要指望编译器呢？通用的规则是，您希望对象的构造函数初始化其每个成员变量，因为未被初始化的数据可能导致难以排除的 `bug`。然而，对于那些将被频繁实例化的小型类，尤其将被实例化为临时对象时，不应为遵守这一规则而牺牲性能。在很多游戏中，这样的类主要是 `Vector` 和 `Matrix`。这些类应提供将元素设置为零的方法（或备用构造函数），而缺省构造函数体应为空。

根据上述原则可以做出这样的推论，即如果额外的构造函数能够提高性能，就应该提供它们。第二个范例中，`Vehicle` 类的代码如下：

```
class Vehicle
{
public:
    Vehicle()
    {
    }

    void SetName(const std::string &name)
    {
        mName = name;
    }

private:
    std::string mName;
};
```

这将创建 `mName`，然后通过 `SetName()` 再次设置其值。同样，与首先创建一个对象，然后调用操作符`=`相比，使用复制构造函数的开销更低。在构造对象时，应这样编写代码——`Vehicle v1(v2)`，而不应这样编写代码——`Vehicle v1; v1=v2;`。

如果要防止编译器自动替您复制对象，则应为相应的类声明一个私有的复制构造函数和操作符`=`，但不实现它们。这样，任何复制该对象的企图都将导致编译错误。另外，还应养成这样的习惯，即显式地声明一个接受单个参数的构造函数，除非您打算将其用于类型转换。这样可以防止编译器在转换类型时生成隐藏的临时对象。

- 预先分配对象并将其保存在缓存中



游戏通常会频繁地分配和释放一些类，如武器或粒子。在 C 游戏中，您通常会预先分配大量的对象，并在必要时使用它们。只要做少量的计划，您就能使用 C++ 完成相同的工作。这里的思想是，您请求一些新的对象，并将旧对象返回到缓存中，而不是不断地创建和销毁对象。可以将缓存实现为一个模板，这样它将适合于任何类，只要这些类有缺省构造函数。本书附带的光盘中包含一个缓存类模板的代码。

您可以在需要对象时进行分配，并将其加入到缓存中，也可以预先分配所有的对象。如果对对象建立了堆栈规则（删除对象  $X$  之前，必须删除在  $X$  对象之后分配的所有对象），则可以在一个连续的内存块中分配缓存。

### 1.1.2 内存管理

与 C 应用程序相比，C++ 应用程序通常更需要注意有关内存管理的细节。在 C 语言中，所有的内存分配和释放工作都是通过 `malloc()` 和 `free()` 显式地完成的，而在 C++ 中，创建临时对象和成员变量时，将隐式地分配内存。大多数 C++ 游戏（和大多数 C 游戏一样）需要自行对内存进行管理。

由于 C++ 游戏很可能进行大量的内存分配，因此应特别注意堆的分段情况。一种选择是采用一种传统方式：在游戏启动后不分配任何内存，或者维护一大块连续的内存，定期地释放它（例如在过关时）。在现代的计算机上，如果您愿意密切注意内存的使用情况的话，没必要采取这种严格的措施。

首先需要覆盖全局操作符 `new` 和 `delete`。通过这两个操作符的自定义实现，使游戏最常用的内存分配工作不用 `malloc()` 完成，而是使用预先分配的内存块。例如，如果您知道每次最多需要分配 10 000 个 4 字节的内存，则应预先分配 40 000 字节的内存，并在必要时使用这些内存块。为跟踪哪些内存块是可用的，可维护一个可用链表（free list），其中每一个可用内存块指向下一个可用的内存块。分配内存时，将第一个内存块从链表中删除；而释放内存时，将其加入到链表的最前面。图 1.1.1 说明了经过分配和释放后，一系列小型内存块是如何形成一个大型内存块的。

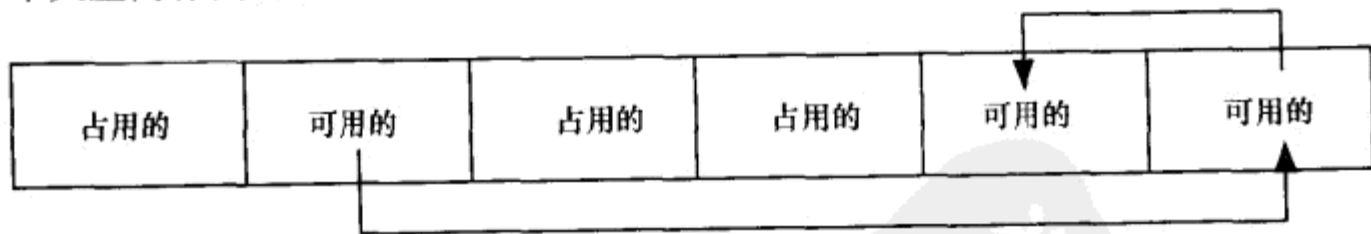


图 1.1.1 可用内存链表

通常您将发现，游戏将分配大量使用时间很短的小型内存块，因此需要为其预留空间。如果预留大量的大型内存块，则当前大量未被使用的内存将被浪费；当这样的内存量达到一定程度后，您将不会使用独立的大型内存块分配函数（`malloc()`）。

### 1.1.3 虚拟函数

反对使用 C++ 编写游戏的人总是指出，虚拟函数是一种神秘的特性，会极大地降低性能。从概念上说，这种机制非常简单。要生成调用对象的虚拟函数，编译器将访问对象的虚拟函数表，从中取得指向该成员函数的指针，建立调用，然后跳到该成员函数所在的地址。在 C 语言的函数调用中，编译器建立调用，然后跳到一个固定的地址，与此相比，虚拟函数调用的额外开销是需要查看虚拟函数表，因为预先并不知道调用的地址。这也将导致处理器指令缓存失的 (cache missing)。

所有重要的 C++ 程序都将使用大量的虚拟函数，因此应防止这些函数对性能带来重大的影响。下面是一个典型的范例：

```
class BaseClass
{
public:
    virtual char *GetPointer() = 0;
};

class Class1 : public BaseClass
{
    virtual char *GetPointer();
};

class Class2 : public BaseClass
{
    virtual char *GetPointer();
};

void Function(BaseClass *pObj)
{
    char *ptr = pObj->GetPointer();
}
```

如果 Function() 对性能有重大影响，则应将虚拟函数 GetPointer 修改为内联函数。为此，方法之一是在 BaseClass 类中增加一个新的受保护的数据成员，内联函数版本的 GetPointer() 返回该数据成员，同时每个类中设置它：

```
class BaseClass
{
public:
    inline char *GetPointerFast()
    {
        return mpPointer;
    }

protected:
    inline void SetPointer(char *pData)
    {
```



```
        mpData = pData;
    }

private:
    char *mpData;
};

// class1 and class2 call SetPointer as necessary
// in member functions

void Function(BaseClass *pObj)
{
    char *ptr = pObj->GetPointerFast();
}
```

另一种动作更大的措施是重新安排类层次结构。如果 Class1 和 Class2 的差别较小，则可以将它们合并为一个类，并使用一个标记来指出您希望该类在运行阶段的行为应该像 Class1 还是 Class2。这样修改后（同时删除纯虚类 BaseClass），也可以将前一个范例中的 GetPointer 作为内联函数。这种变换并不好，但当它在循环体内且在缓存较小的计算机上运行时，您将会愿意这样做，已避免调用虚拟函数。

虽然新增一个虚拟函数只会增加指向类表的指针的大小（通常开销可忽略不计），但类中的第一个虚拟函数要求每个对象独有一个指向虚拟函数表的指针。这意味着对于常用的小型类不能有任何虚拟函数，因为在这种类中，这样的开销是不可接受的。由于继承通常要求使用一个或多个虚拟函数（至少有一个虚拟析构函数），因此对于被大量使用的小型对象，不能继承。

#### 1.1.4 代码长度

C++编译器生成的代码过大。由于内存有限，且代码越小速度越快，因此使可执行文件尽可能小显得至关重要。首先要调整好编译器。如果编译器将调试信息存储在可执行文件中，则应关闭生成调试信息的功能（Microsoft Visual C++将调试信息存储在独立于可执行文件的文件中，因此没必要这样做）。异常处理也会生成额外的代码，因此应尽可能地除掉异常处理生成的代码。应确保将链接程序配置成将未被使用的函数和类删除。启用编译器的最高级别优化，并将其设置为对代码长度而不是速度进行优化——有时候这将实际使得生成的代码速度更快，因为指令缓存的相关性更强（如果使用这种设置，一定要确保本征函数也被启用）。请删除调试打印语句中浪费空间的字符串，并让编译器将重复的常量字符串合并为单个实例。

内联技术常常是导致大型函数的罪魁祸首。编译器可能遵守或忽略关键字 `inline`，也可能将函数作为内联函数而不告知您。这也是应使构造函数尽可能小的另一个原因，这可以防止堆栈上的对象将生成大量的内联代码。另外，还应注意重载操作符，如果 `m2` 和 `m3` 是矩阵，则诸如 `m1 = m2 * m3` 这样简单的表达式将生成大量的内联代码。因此，一定要详细了解编译器中关于内联函数的设置。

启用运行阶段运行信息（RTTI），需要编译器为程序中的每个类生成一些静态信息。启用 RTTI 通常是为了让代码能够调用 `dynamic_cast`，并确定对象的类型。应完全避免 RTTI 和

`dynamic_cast`，以节省空间（另外，在有些实现中，`dynamic_cast` 的开销非常大）。当您确实需要根据类型而采取不同的行为时，应添加一个行为因类型而不同的虚拟函数，这是一种更好的面向对象的设计方式（对于 `static_cast`，这一点不适用，从性能上说，这种转换同 C 风格的类型转换相同）。

### 1.1.5 标准模板库

标准模板库（STL）是一组模板，它们实现了常用的数据结构和算法，如动态数组（被称为矢量）、集合和映射表（`map`）。使用 STL 可以节省大量的时间，避免自己编写和调试这些容器。同样，如果您希望获得最高的效率，则需要了解有关 STL 实现的细节。

为使实现的适用范围尽可能大，STL 标准没有在内存分配方面做任何规定。STL 容器的所有操作都有一定的性能保证，例如，集合的插入操作的时间为  $O(\log n)$ ，但容器的内存使用情况没有任何保证。

下面详细介绍游戏开发中一个常见的问题：您想存储大量的对象（我们将其称为一串对象，虽然不一定要将其存储在一个 STL 链表中）。通常，您希望每个对象只在链表中出现一次，这样，如果该对象已经在集合中，则无须担心无意间再次将其插入。STL 集合将忽略重复的元素，其插入、删除和查找的速度为  $O(\log n)$ ，因此是最佳选择，是吗？

也许是。虽然集合的大部分操作的速度为  $O(\log n)$ ，但这种表示法隐藏了常量可能非常大的情况。虽然集合使用内存的方式取决于实现，但很多实现都是基于红-黑树（`red-black tree`）的，其中每个树节点存储集合中的一个元素。通常，每次插入一个元素时，都分配一个树节点，而每当元素被删除时，都释放一个树节点。根据插入和删除元素的频率，内存分配函数所花的时间可能比使用集合时算法节省的时间要多。

另一种解决方案是使用 STL 矢量来存储元素。矢量确保在集合的最后插入元素的时间是固定的。这意味着矢量通常只是偶尔重新分配内存，例如被填满后将长度翻倍。使用矢量来存储一系列不重复的元素时，您首先对矢量进行检查，看元素是否在其中，如果不在，则将元素添加到最后面。检查整个矢量所需的时间为  $O(n)$ ，但其涉及的常量可能很小。由于矢量的所有元素通常被连续存储在内存中，因此，检查整个矢量是一个缓存友好（`cache-friendly`）的操作。检查整个集合时，可能需要不断地在内存中移动，因为红-黑树的各个元素可能分散在内存的各个地方。另外，集合建立树的开销也非常大。如果要存储的全部是对象指针，则使用集合时，需要的内存量和可能是使用矢量时的 3~4 倍。

从集合中删除一个元素的时间为  $O(\log n)$ ，在不考虑对 `free()` 的调用时，这看起来很快。从矢量中删除一个元素的时间为  $O(n)$ ，因为从该元素开始到矢量最后的每个元素都将被复制到前一个位置。然而，如果矢量中的所有元素都是指针，则整个复制工作只需调用一次 `memcpy()` 便可完成，这通常非常快（这也是通常应将指向对象的指针，而不是对象本身存储在 STL 集合中的原因之一。如果直接存储对象本身，在执行诸如复制等操作时，将调用大量的构造函数）。

如果您还不相信集合和映射表带来的麻烦比好处多，请看遍历集合的开销：

```
for (Collection::iterator it = collection.begin();
     it != collection.end(); ++it)
```

如果 Collection 是矢量，则++it 将是指针递增操作——一条机器指令。但如果 Collection 为集合或映射表，则++it 将需要遍历到红-黑树的下一个节点，这是一个比较复杂的操作，同时更有可能导致缓存失的，因为树节点可能分散在整个内存中。

当然，如果需要将大量的元素存储在集合中，并需要执行大量的成员查询操作，则与内存方面的开销相比，集合的  $O(\log n)$  性能带来的好处更大；同样，如果只是偶尔使用集合，则性能方面的差异将是无关紧要的。您应进行性能方面的测试，确定当  $n$  为多大时，集合的速度将更快。您可能惊讶地发现，对于您的游戏通常使用的所有  $n$  值而言，使用矢量的性能都优于使用集合。

然而，有关 STL 内存使用方面的内容并不止这些。知道调用 clear() 方法时，集合是否释放其使用的内存至关重要。如果不释放，将产生内存碎片。例如，如果游戏启动时，矢量为空，随着游戏的进行，元素被加入到矢量中，当玩家重新启动游戏时，将调用 clear() 方法，而矢量可能根本不释放其使用的内存。空矢量可能仍然占据堆中的某些空间，从而产生碎片。对于这种问题，解决方式有两种。首先，可以在矢量被创建时调用 reserve()，预留足够多的空间，以存储可能需要的最大数量的元素。如果这种方式不可行，则可以采用下面的方式，显式地强制矢量释放其使用的内存：

```
vector<int> v;  
// ... elements are inserted into v here  
vector<int>().swap(v); // causes v to free its memory
```

集合、链表和映射表通常不存在这样的问题，因为它们为每个元素分别分配和释放内存。

### 1.1.6 高级特性

语言拥有某种特性，并不意味着您一定要使用该特性。有些看起来非常简单的特性其性能可能非常差，而看起来非常复杂的特性其实际性能却可能非常好。C++ 最大的问题在于它严重依赖于编译器，因此使用编译器之前，一定要了解其开销。

C++ 字符串便是一个这样的例子，从理论上说，这种特性非常棒，但如果性能至关重要时，应避免使用。请看下面的代码：

```
void Function(const std::string &str)  
{  
}  
  
Function("hello");
```

调用 Function() 时，由于提供的参数是一个 const char \*，将调用字符串的构造函数。在一种商用实现中，该构造函数执行 malloc()、strlen() 和 memcpy() 函数，而析构函数立刻执行一些非常重要的工作（因为该实现的字符串是引用计数的），然后执行 free()。分配的内存完全是浪费，因为程序的数据段中已经包含字符串“hello”，已在内存中有效地将其进行了复制。如果 Function() 被声明为带有一个 const char \* 参数，则该调用将不会带来这些开销。为方便操纵字符串，付出的代价是非常高昂的。

模板体现了效率的两个极端。根据语言标准，当使用特定的类型实例化模板时，编译器将生成模板的代码。从理论上说，好像一个模板声明将导致大量几乎相同的代码。如果有一个由指向 Class1 的指针组成的矢量和一个由指向 Class2 的指针组成的矢量，则最终的可执行文件中将包含 `vector` 的两个拷贝。

实际上，在大部分编译器中，情况要好些。首先，仅当模板成员函数被调用后，才会生成其代码。其次，如果维持正确的行为，编译器可以只生成代码的一个拷贝。您通常将发现，在前面有关矢量的范例中，将只会生成代码的一个拷贝（可能是 `vector<void*>` 的代码）。使用优秀的编译器时，模板可以提供通用编程的所有好处，同时确保较高的性能。

C++ 的一些特性（如初始化列表和先递增）通常能够提高性能，而其他特性（如操作符重载和 RTTI）看起来好像没有什么问题，但将严重地降低性能。STL 集合说明了这样一点，即盲目地信任函数的算法运行时间，将使您误入歧途。应避免使用语言和库中速度可能很慢的特性，并花一些时间来熟悉编译器和配置文件（`profiler`）中的选项，您将很快学会如何提高速度并找到游戏中性能方面的问题。

### 1.1.7 参考文献

Cormen, Thomas, Charles Leiserson, and Ronald Rivest, *Introduction to Algorithms*, Cambridge, Massachusetts, MIT Press, 1990.

Isensee, Peter, C++ Optimization Strategies and Techniques, [www.tantalon.com/pete/cppopt/main.htm](http://www.tantalon.com/pete/cppopt/main.htm).

Koenig, Andrew, "Pre- or Postfix Increment," *The C++ Report*, June, 1999.

Meyers, Scott, *Effective C++*, Second Edition, Reading, Massachusetts: Addison-Wesley Publishing Co., 1998.

Sutter, Herb, Guru of the Week #54: Using Vector and Deque, [www.gotw.ca/gotw/054.htm](http://www.gotw.ca/gotw/054.htm).

感谢 Pete Isensee 和 Christopher kirmse 审阅本文。



## 1.2 内联函数和宏

Peter Dalton, Evans & Sutherland

pdalton@xmission.com

在游戏编程中，无论怎样强调函数的速度和效率都不过分，尤其是那些在每一帧中需要执行多次的函数更是如此。许多程序员在处理常用的、时间—关键的例程时大量地使用宏，因为宏消除了函数所需的调用/返回序列，而这些序列对函数调用的开销而言是敏感的。然而，使用编译指令`#define`来实现类似于函数的宏更容易引起问题，不值得这样做。

### 1.2.1 内联函数的优点

通过使用内联函数，可以轻易地消除宏固有的许多缺陷。以下面的宏定义为例：

```
#define max(a,b) ((a) > (b) ? (a) : (b))
```

如果这样调用宏：`max(++x, y)`，将发生什么情况呢？如果  $x=5$ ， $y=3$ ，则该宏将返回 7，而不是期望的 6。这说明了宏最常见的副作用——作为实参传递的表达式可能被计算多次。要避免这种问题，可以使用内联函数来实现上述目标：

```
inline int max(int a, int b) { return (a > b ? a : b); }
```

通过使用内联函数，可以确保所有的参数都只被计算一次，因为根据定义，参数必须遵守强加给普通函数的所有规则和类型安全性。

宏存在的另一个问题是操作符的优先级，下面的宏说明了这一点：

```
#define square(x) (x*x)
```

如果使用表达式 `2+1` 调用上述宏，返回的结果将为 5，而不是期望的 9。这里的问题是，乘法运算符的优先级高于加法运算符。虽然使用括号将整个表达式括起可以解决这个问题，但使用内联函数更容易避免这样的问题。

宏的另一个缺陷是需要处理多条语句的宏，并确保宏中的所有语句都被正确地执行。来看下面这个简单的宏，该宏用于将任何数字转换为 0~1 之间的值：

```
#define clamp(a) \
    if (a > 1.0) a = 1.0; \
```



```
if (a < 0.0) a = 0.0;
```

如果将上述宏用于下面的 for 循环中：

```
for (int ii = 0; ii < N; ++ii)
    clamp( numbersToBeClamped[ii] );
```

则如果数字小于 0，则不会被转换。仅当 `ii==N`，导致 for 循环结束后，表达式 `if(numbersToBeClamped[i]<0.0)` 才会被计算。这也会引起问题，因为此时索引变量已经超出范围，很可能导致内存越界错误，进而导致程序崩溃。虽然使用内联函数替代宏来完成同样的功能不是惟一的解决方案，但这种方案最为整洁。

鉴于宏存在这些固有的缺陷，我们来看看内联函数的优点：

- 内联函数遵守强加给函数的所有规则和类型安全性，这可以避免将意外或非法的参数作为实参来传递它。
- 定义内联函数的语法和其他函数相同，只是在函数声明中使用了关键字 `inline`。
- 将表达式作为实参传递给内联函数时，将在进入函数体之前被计算，因此只被计算一次。正如前面指出的，传递给宏的表达式可能被计算多次，从而导致不安全和意外的副作用。
- 可以使用诸如微软的 Visual C++ 等调试器对内联函数进行调试；而这对于宏而言是不可能的，因为宏将在分析程序接手前和程序符号创建前被扩展。
- 内联函数可以提高程序的可读性和可维护性，因为它们使用的语法和常规函数调用相同，而且不会意外地修改参数。

另外，内联函数消除了函数调用的开销，因此性能比常规函数高。这些开销包括栈结构的建立、参数传递、栈结构的恢复以及返回等任务。除了上述重要的优点外，内联函数还使编译器能够对代码进行优化。通过将内联函数替换为代码，插入的代码可以被进一步地优化，而这本来是不可能的，因为大多数编译器不会执行过程间的优化。让编译器进行全局优化，如消除相同的子表达式、删除循环中不变的量，可以极大地提高性能，缩短代码的长度。

对内联函数的限制中，惟一在宏中不存在的是参数的类型。宏可以接受任何类型的参数，而内联函数只能接受指定类型的参数，以实现类型安全性。可以通过使用内联模板函数来克服这种局限性，内联模板函数可以接受任何类型的参数，实现类型安全性，同时具备内联函数的所有优点。

### 1.2.2 何时使用内联函数

为何不将所有的函数都作为内联函数呢？这可以消除整个程序的函数开销，从而使填充速度更快，响应时间更短吗？显然回答是否定的。虽然通过消除函数开销以及允许执行过程间编译器优化，代码扩展可以提高速度，但这是以增加代码长度为代价的。度量程序的性能时，需要考虑两个因素：执行速度和代码的实际长度。代码长度增加将占用更多的宝贵内存，从而降低执行速度。随着程序的内存需求增加，缓存失的（cache missing）和缺页的可能性也将增大。缓存失的将导致些微的延迟，而缺页将导致很大的延迟，因为虚拟内存单元不在物理内存中，必须从硬盘中取回它。在 400 MHz 的 Pentium II 台式机上，一次硬缺页将导致



大约 10 $\mu$ s 的延迟，这大约相当于 4 000 000 个 CPU 周期[Heller99]。

如果说内联函数并非任何时候都是合适的，那么何时应使用它们呢？这个问题的答案取决于具体情况，必须由程序员来作出判断。然而，下面是一些何时应使用内联函数的指导原则：

- 小型方法，如私有数据成员的存取函数（accessor）；
- 返回对象状态信息的函数；
- 小型函数，通常不超过三行；
- 不断被调用的小型函数，例如，在时间关键的渲染循环中。

对于较长的函数，其调用/返回序列所用时间的比例较低，将其作为内联函数带来的好处较小。然而，在使用正确的情况下，内联将极大地提高程序的性能。

### 1.2.3 何时使用宏

---

虽然宏存在一些问题，但在为数不多的几种情况下，它们是非常有用的。例如，可以使用宏来创建小型伪语言，这种伪语言的功能可以非常强大。一组宏能够提供这样的框架，即使得创建状态机变成小菜一碟，同时具备非常高的可调试性和健壮性。有关这种技术的优秀范例请参阅本文最后的文献[Rabin00]；另一个范例是在屏幕上打印枚举类型，如下所示：

```
#define CaseEnum(a)    case(a) : PrintEnum( #a )
switch (msg_passed_in) {
    CaseEnum( MSG_YouWereHit );
        ReactToHit();
        break;
    CaseEnum( MSG_GameReset );
        ResetGameLogic();
        break;
}
```

其中，PrintEnum()是一个宏，它在屏幕上打印一个字符串。#是一个字符串化操作符，将宏的参数转换为字符串常量[MSDN]。因此，无需创建查找表（其中包含枚举数到字符串的映射，通常难以维护），以检索宝贵的调试信息。

避开宏存在的问题的关键是：首先要了解问题；其次是知道其他的实现方式。

### 1.2.4 微软特有的情况

---

除了标准的 inline 关键字外，Microsoft's Visual C++编译器还支持另外两个关键字。关键字 \_\_inline 命令编译器进行成本/收益分析，仅当有益时，才将其作为内联函数。关键字 \_\_forceinline 命令编译器总是将函数作为内联函数。即使使用这些关键字，在有些情况下编译器也可能不遵守，微软公司的文档[MSDN]对此进行了说明。

### 1.2.5 参考文献

---

[Heller99] Heller, Martin, *Developing Optimized Code with Microsoft Visual C++ 6.0*, Microsoft MSDN Library, January 2000.

[McConnell93] McConnell, Steve, *Code Complete*, Microsoft Press, 1993.

[MSDN] Microsoft Developer Network Library, <http://msdn.microsoft.com>.

[Myers98] Myers, Scott, *Effective C++, Second Edition*, Addison-Wesley Longman, Inc., 1998.

[Rabin00] Rabin, Steve, "Designing a General Robust AI Engine," *Game Programming Gems*. Charles River Media, 2000; pp. 221~236.



## 1.3 抽象接口编程

Noel Liopis, Meyer/Glass Interactive  
nllopis@mgigames.com

**抽**象接口的概念很简单，但功能非常强大。它让我们能够将接口完全同其实现分开，这带来了一些好处：

- 易于在代码的不同实现间切换，而不影响游戏的其他部分。当试验不同的算法或在不同的平台上修改实现时，这尤其有用。
- 可以在运行阶段修改实现。例如，如果图形渲染器（renderer）是通过抽象接口实现的，则可以在游戏运行期间，在软件渲染器和硬件加速渲染器之间进行切换。
- 对接口的用户完全隐藏了实现的细节。这样在整个工程中包含的头文件将更少，重新编译的速度将更快，因此当需要完全重新编译整个工程时，所需的时间将更短。
- 即使在游戏被编译和发行之时，还可以很容易地将已有接口的新实现加入到游戏中。这样，可以通过提供更新或用户自定义的修改，轻松地对游戏进行扩展。

### 1.3.1 抽象接口

在 C++ 中，抽象接口是一个只有公有纯虚函数的基类。纯虚函数是一种没有实现的虚拟成员函数。派生类必须实现这些函数，否则编译器将禁止对该类进行实例化。纯虚函数在声明后面被加有=0。

下面是最简单的音响系统的一个抽象接口范例。该接口将在一个头文件中自行声明：

```
// In SoundSystem.h
class ISoundSystem {
public:
    virtual ~ISoundSystem() {};
    virtual bool PlaySound ( handle hSound ) = 0;
    virtual bool StopSound ( handle hSound ) = 0;
};
```

抽象接口没有提供任何实现，而只是定义了使用该音响系统时必须遵守的规则。接口的用户只要了解 ISoundSystem，便可以使用任何音响系统实现。

下面的头文件指出了上述接口的一个实现范例：

```
// In SoundSystemSoftware.h
#include "SoundSystem.h"

class SoundSystemSoftware : public ISoundSystem {
public:
    virtual ~SoundSystemSoftware();
    virtual bool PlaySound ( handle hSound );
    virtual bool StopSound ( handle hSound );

    // The rest of the functions in the implementation
};
```

显然，对于 `corresponding.cpp` 文件中的每个函数，都必须提供实际的实现。要使用这个类，必须这样做：

```
ISoundSystem * pSoundSystem = new SoundSystemSoftware();
// Now we're ready to use it
pSoundSystem->PlaySound ( hSound );
```

那么，以这种迂回的方式创建音响系统后，完成了哪些工作呢？几乎完成了本文开始所指出的一切。

- 很容易创建该音响系统的另一种实现（可能是一个硬件版本），所需要做的只是从 `ISoundSystem` 派生出一个新类，并实例化它。其他的一切将同原来相同，而无需做更多的修改。
- 在运行阶段，可以在两个类之间切换。只要 `pSoundSystem` 指向的是一个有效的对象，程序的其他部分将不知道使用的是哪个类，因此我们可以随心所欲地修改它们。显然，我们必须注意具体类的限制。例如，有些类将保存一些状态信息，或者要求首次被使用之前进行初始化。
- 我们对用户隐藏了所有的实现细节。通过实现接口，我们承诺提供文档指出的行为，而不管实现是什么样的。与使用 `if` 语句来检查音响系统的类型相比，代码更为清晰；维护代码的工作也容易得多。

### 1.3.2 添加一个工厂 (factory)

还有一个细节没有介绍：我们还未对用户完全隐藏具体的实现。毕竟，用户仍将对他们要使用的类的具体实现执行 `new` 操作。这里的问题是，用户需要使用 `#include` 将包含实现声明的头文件包含进去。不幸的是，由于 C++ 的设计，当用户包含头文件时，也将获得许多关于类实现细节的信息，而本来他们不应知道这些信息。用户将看到所有的私有和保护成员，甚至可能需要包含进只在类实现中使用的头文件。

更糟的是，接口用户知道接口指针指向的类的类型，从而可能试图将其转换为实际的类型以获取一些“特殊的特性”，或者依赖于实现特有的行为。在这样的情况下，将丧失使用抽象接口的好处，因此应尽可能避免这样的情况发现。

对此，解决之道是使用一个抽象工厂（abstract factory）[Gamma95]。抽象工厂是一个类，其惟一的用途是实例化接口的具体实现。下面是前面的音响系统的一个基本工厂范例：

```
// In SoundSystemFactory.h
class ISoundSystem;

class SoundSystemFactory {
public:
    enum SoundSystemType {
        SOUND_SOFTWARE,
        SOUND_HARDWARE,
        SOUND_SOMETHINGELSE
    };

    static ISoundSystem * CreateSoundSystem(SoundSystemType type);
};

//In SoundSystemFactory.cpp
#include "SoundSystemSoftware.h"
#include "SoundSystemHardware.h"
#include "SoundSystemSomethingElse.h"

ISoundSystem * SoundSystemFactory::CreateSoundSystem ( SoundSystemType_type )
{
    ISoundSystem * pSystem;

    switch ( type ) {
        case SOUND_SOFTWARE:
            pSystem = new SoundSystemSoftware();
            break;
        case SOUND_HARDWARE:
            pSystem = new SoundSystemHardware();
            break;
        case SOUND_SOMETHINGELSE:
            pSystem = new SoundSystemSomethingElse();
            break;
        default:
            pSystem = NULL;
    }

    return pSystem;
}
```

这样，问题便解决了。用户只需包含头文件 SoundSystemFactory.h 和 SoundSystem.h。事实上，我们甚至无需使其他头文件可用。要使用具体的音响系统，用户可以这样编写代码：

```
ISoundSystem * pSoundSystem;  
pSoundSystem = SoundSystemFactory::CreateSoundSystem  
                (SoundSystemFactory::SOUND_SOFTWARE);  
// Now we're ready to use it  
pSoundSystem->PlaySound ( hSound );
```



抽象接口中一定要包含一个虚拟析构函数。如果没有，C++将自动生成一个非虚拟析构函数，这将导致具体实现的析构函数不被调用（这将是一个 bug，并且难以发现）。和常规的成员函数不同，我们不能只提供一个纯虚析构函数，因此需要创建一个空函数，以满足编译器的需求。

### 1.3.3 抽象接口特性

另一种看待抽象接口的方式是，将其视为一组行为。类实现接口时，是承诺以特定的方式运行的。例如，下面的接口将由能够被渲染到屏幕上的对象使用：

```
class IRenderable {  
public:  
    virtual ~IRenderable() {};  
    virtual bool Render () = 0;  
};
```

我们可以设计一个类来表示这样的 3D 对象，即从 `IRenderable` 派生而来，且提供了将自己渲染到屏幕上的方法。同样，我们还可以有地形（`terrain`）类，它也是从 `IRenderable` 派生而来的，但提供的渲染方法完全不同。

```
class Generic3DObject : public IRenderable {  
public:  
    virtual ~Generic3DObject();  
    virtual bool Render();  
  
    // Rest of the functions here  
};
```

渲染循环将遍历所有的对象，如果被遍历的对象能够被渲染，则调用其 `Render()` 方法。同样，这种接口的强大功能也是来自对接口隐藏真正的实现：现在，可以新增一种全新的对象，只要它表现为 `IRenderable` 接口，渲染循环便可以像渲染其他对象那样渲染它。如果不使用抽象接口，渲染循环将知道具体的对象类型（普通的 3D 对象、地形还是其他），以决定是否调用其特有的渲染函数。在这种情况下，创建一种新的、具备渲染功能的对象时，将需要修改渲染循环以及代码的许多其他部分。

我们可以通过检查对象是否是从 `IRenderable` 派生而来的，以确定它是否可被渲染。不幸的是，这要求在编译代码时，将编译器的 RTTI（运行阶段类型信息）选项打开。启用 RTTI 通常会降低性能，增加所需的内存，因此很多游戏将这项功能关闭。可以使用自定义的 RTTI，不过我们还是采用 COM（组件对象模型）的方式，并提供一个 `QueryInterface` 函数 [Rogerson97]。



如果对象实现了特定的接口，QueryInterface 将把输入指针转换为这种接口，并返回 true。要创建自己的 QueryInterface 函数，需要一个基类，所有继承一组接口的相关对象都将从该基类派生而来。我们甚至可以将基类本身编写为一个类似于 COM 的 Iunknown 的接口，但如果这样做将使问题复杂化。

```
class GameObject {
public:
    enum GameInterfaceType {
        IRENDERABLE,
        IOTHERINTERFACE
    };

    virtual bool QueryInterface (const GameInterfaceType type,
                                void ** pObj );

    // The rest of the GameObject declaration
};
```

对于简单的游戏对象而言，QueryInterface 的实现没有什么价值。因为它没有实现任何接口，且总是返回 false。

```
bool GameObject::QueryInterface (const GameInterfaceType type,
                                void ** pObj ) {
    return false;
}
```

3D 对象类的实现不同于 GameObject，因为它将实现 Irenderable 接口。

```
class 3DObject : public GameObject, public IRenderable {
public:
    virtual ~3DObject();
    virtual bool QueryInterface (const GameInterfaceType type,
                                void ** pObj );

    virtual bool Render();
    // Some more functions if needed
};

bool 3DObject::QueryInterface (const GameInterfaceType type,
                                void ** pObj ) {
    bool bSuccess = false;
    if ( type == GameObject::IRENDERABLE ) {
        *pObj = static_cast<IRenderable *>(this);
        bSuccess = true;
    }
    return bSuccess;
}
```

3DObject 类必须负责覆盖 QueryInterface，检查它支持哪些接口，并执行相应的类型转换。现在，我们来看看渲染循环，它简单、灵活，且对要渲染的对象的类型一无所知。

```
IRenderable * pRenderable;
```

```
for ( all the objects we want to render ) {  
    if ( pGameObject->QueryInterface (GameObject::IRENDERABLE,  
        (void**)&pRenderable) )  
    {  
        pRenderable->Render();  
    }  
}
```

现在，便可以实现本文开始列出的抽象接口的最后一条承诺了：轻松地添加新的实现。对于这样的渲染循环，如果我们给它提供新的对象类型，且它们实现了 **IRenderable** 接口，则一切都将按预期的进行，而无需对渲染循环进行修改。引入新的对象类型的最简单的方式是，使用更新后的库或包含新类的代码，重新链接工程。我们将通过 DLL（或目标平台中相对应的机制），在运行阶段引入新的对象类型（这已经超过本文所讲述的范围）。这种改进将让我们能够发布新的游戏对象或游戏更新，而无需对可执行文件进行修补。用户则可以使用这种方式，轻松地对游戏进行修改。

请注意，我们完全可以继承多个接口。这样做意味着继承多个接口的类将提供每个接口指定的所有服务。例如，对于要进行碰撞检测（collision detection）的对象，我们可以提供一个 **ICollidable** 接口。3D 对象可以同时继承 **IRenderable** 和 **ICollidable**，但表示烟雾的类则只继承 **IRenderable**。

然而，需要注意的是：虽然使用多个抽象接口的技术功能强大，但也可能导致设计过度复杂，从而相对于使用单继承的设计，没有任何优势可言。另外，对于动态特征来说，多重继承不太适合，因而只能用于对象固有的永久性特征。

虽然很多人建议远离多重继承，但还是要看它是否有用，是否没有重大缺陷。最多从一个真正的父类和多个接口函数进行继承，并不会导致令人恐惧的金字塔形继承树（即双亲的父节点均为同一个类）或多重继承的其他许多常见的缺陷。

### 1.3.4 一切都是代价的

至此，我们介绍了抽象接口许多吸引人的特性。然而，所有这些特性都是需要付出代价的。大多数时候，使用抽象接口带来的好处超过了其潜在的问题，但了解这种技术的缺陷和局限性至关重要。

首先，设计将更复杂。对于不熟悉抽象接口的人来说，首次看到多余的类以及对接口的查询可能会感到迷惑。仅当它能够带来截然不同的效果时，才应采用，而不能不分青红皂白，在整个游戏中都进行使用；否则将使设计晦涩难懂，并带来破坏。

通过使用抽象接口，可以隐藏所有的私有实现，但这将使它们更难调试。如果拥有的只是一个 **IRenderable \***类型的变量，则不经过大量繁琐的类型转换，将无法在调试器的交互式监视窗口中看到该变量指向的实际对象的私有内容。另一方面，大多数时候我们无须担心这一点，因为实现已经被隔离开并经过仔细的测试，我们只需正确地使用接口即可。

另一个缺点是，无法通过继承对已有的接口进行扩展。回到第一个范例，我们可能想对 **SoundSystemHardware** 类进行扩展，以添加几个该游戏特有的函数。不幸的是，我们无法访问这个类的实现，当然也就无法继承或扩展它。但我们仍然可以通过使用一个派生类来修改

已有的接口或提供一个新的接口，但这需要在实现中去完成，而不能在游戏代码中完成。

最后，请注意，抽象接口中的每个函数都是虚拟函数。这意味着每当计算机通过抽象接口调用这些函数时，间接程度都将增加一级。对于现代的计算机和游戏控制台而言，只要对那些在内嵌循环中调用的函数不使用接口，这通常不是什么问题。例如，创建一个包含函数 `DrawPolygon()` 或 `SetScreenPoint()` 的接口可能不是一个好主意。

### 1.3.5 结论

---

抽象接口是一种功能强大的技术，我们可以大量使用它，而不会带来太多的开销，也无需进行结构性修改。知道如何正确地使用它以及在什么时候应采用其他方式至关重要。抽象接口最适合用于可替换的模块（图形渲染器、空间数据库、AI 行为）、各种可插入的模块或用户可扩展的模块（工具扩展、游戏行为）。

### 1.3.6 参考文献

---

- [Gamma95] Gamma, Eric et al, *Design Patterns*, Addison-Wesley, 1995.
- [Lakos96] Lakos, John, *Large Scale C++ Software Design*, Addison-Wesley, 1996.
- [Rogerson97] Rogerson, Dale, *Inside COM*. Microsoft Press, 1997.



## 1.4 从 DLL 中导出 C++ 类

Herb Marselas, Ensemble Studios

hmarselas@ensemblestudios.com

从动态链接库（DLL）中导出 C++ 类供其他应用程序使用，是一种封装实例化（instanced）功能或共享可派生（derivable）功能，而又无需共享导出类的源代码的简单方式。这种方式与 Microsoft COM 有些类似，但量级更轻，派生起来更容易，提供的接口也更简单。

### 1.4.1 导出函数

从本质上说，从 DLL 中导出函数和导出类之间并没有什么区别。为从一个 DLL 中导出函数 myExportedFunction，需要在 DLL 工程的预处理器选项中定义值 BUILDING\_MY\_DLL，而不是在使用该 DLL 的工程中定义。这将导致联编该 DLL 时，DLLFUNCTION 将被替换为 \_\_declspec(dllexport)，而联编使用该 DLL 的工程时，将被替换为 \_\_declspec(dllimport)。

```
#ifdef BUILDING_MY_DLL
#define DLLFUNCTION __declspec(dllexport) // defined if building the
                                           // DLL
#else
#define DLLFUNCTION __declspec(dllimport) // defined if building the
                                           // application
#endif
DLLFUNCTION long myExportedFunction(void);
```

### 1.4.2 导出类

从 DLL 中导出 C++ 类稍微复杂些，因为存在多种方式。最简单的情况是，类本身被导出。和前面一样，DLLFUNCTION 宏用来声明被 DLL 导出（或被应用程序导入）的类。

```
#ifdef BUILDING_MY_DLL
#define DLLFUNCTION __declspec(dllexport)
#else
#define DLLFUNCTION __declspec(dllimport)
#endif
```

```

class DLLFUNCTION CMyExportedClass
{
public:
    CMyExportedClass(void) : mdwValue(0) { }

    void setValue(long dwValue) { mdwValue = dwValue; }
    long getValue(void) { return mdwValue; }
    long clearValue(void);

private:
    long mdwValue;
};

```

如果包含这个类的 DLL 被隐式地链接(换句话说,工程是使用该 DLL 的 lib 文件链接的),则使用该类时只需声明其一个实例 CMyExportedClass 即可。这也使得继承该类时,就像它是在应用程序中直接声明的一样。在应用程序中声明派生类时,通常无需做其他的声明:

```

class CMyApplicationClass : public CMyExportedClass
{
public:
    CMyApplicationClass(void)
    {
    }
    ...
};

```

在应用程序中声明或分配一个从 DLL 中导出的类存在一个问题:可能使内存跟踪程序迷惑,使其错误地报告内存的分配或释放情况。为解决这种问题,必须在 DLL 中添加分配和释放被导出类的实例的助手(helper)函数。导出类的用户必须调用分配函数来创建这个类的实例,并调用删除函数来释放实例。当然,这样做的缺点是,无法在应用程序中通过导出类派生出其他的类。如果从导出类派生出其他类很重要,且该工程使用了内存跟踪程序,则或者该内存跟踪程序必须知道发生的情况,或者使用新的内存跟踪程序替换它。

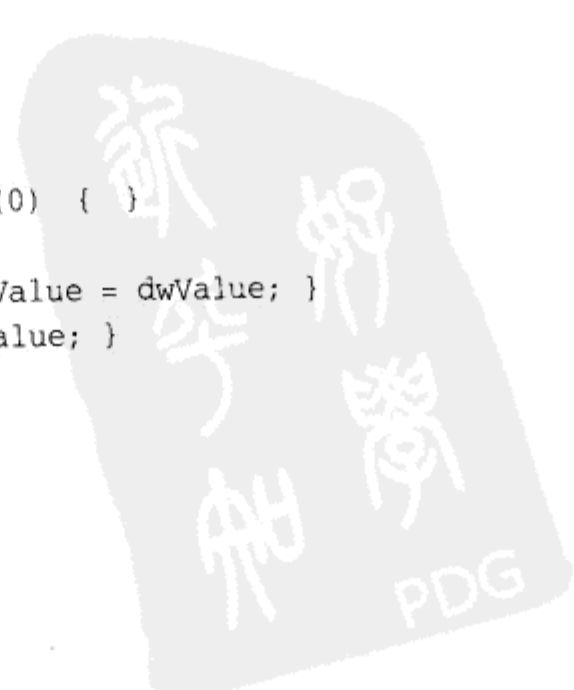
```

#ifdef _BUILDING_MY_DLL
#define DLLFUNCTION __declspec(dllexport)
#else
#define DLLFUNCTION __declspec(dllimport)
#endif

class DLLFUNCTION CMyExportedClass
{
public:
    CMyExportedClass(void) : mdwValue(0) { }

    void setValue(long dwValue) { mdwValue = dwValue; }
    long getValue(void) { return mdwValue; }
    long clearValue(void);

```



```
private:
    long mdwValue;
};

CMyExportedClass *createMyExportedClass(void) {
    return new CMyExportedClass; }
void deleteMyExportedClass(CMyExportedClass *pclass) {
    delete pclass; }
```

### 1.4.3 导出类成员函数

即使添加助手函数后，用户仍然不用调用助手函数 `createMyExportedClass` 就可以创建类的实例，因为类本身被从 DLL 中导出。这种问题很容易解决，只需将导出规范从类级移到用户需要访问的各个函数即可。这样，使用该类的应用程序将无法创建类本身的实例，而必须调用助手函数 `createMyExportedClass` 来创建类的实例，并调用 `deleteMyExportedClass` 来释放类实例。

```
class CMyExportedClass
{
public:
    CMyExportedClass(void) : mdwValue(0) { }

    DLLFUNCTION void setValue(long dwValue) { mdwValue = dwValue; }
    DLLFUNCTION long getValue(void) { return mdwValue; }
    long clearValue(void);

private:
    long mdwValue;
};

CMyExportedClass *createMyExportedClass(void) {
    return new CMyExportedClass; }
void deleteMyExportedClass(CMyExportedClass *pclass) {
    delete pclass; }
```

还应该注意的是，虽然 `CMyExportedClass::clearValue` 是一个公有成员函数，但类用户在 DLL 的外面不能调用它，因为它没有被声明为 `dllexport`。这对于那些需要使其一些函数对 DLL 外的类用户是公用的，同时使其他公有函数只能在 DLL 内部可用的复杂类而言，这是一个功能强大的工具。采用这种策略的一个范例是 SDK for Discreet's 3D Studio MAX。大多数类都有被导出和未被导出的函数，这使得用户能够访问被导出的成员函数或从其派生出其他功能，同时让 SDK 的开发人员拥有一组只有在内部才可以使用的成员函数。

### 1.4.4 导出虚拟类成员函数

对于使用 Microsoft Visual C++ 6 的用户来说，应注意一个潜在的问题。如果您试图导出



类的成员函数，而没有使用导出该类的 DLL 的 lib 文件来链接（在运行阶段使用 LoadLibrary 来装载该 DLL），则如果没有启用内联函数扩展功能，将对于您引用的每个函数，都将出现一个“unresolved external symbol”的错误。无论函数是否在头文件中被声明，都会出现这样的情况。对此，一种解决方法是，将内联函数扩展选项修改为“Only \_\_inline”或“Any Suitable”。不幸的是，这可能与您希望在调试联编中关闭内联函数扩展相冲突。另一种解决方案是，将函数声明为虚拟的，这样将生成正确的代码，而不管内联函数选项的设置是什么。在很多情况下，您可能想将导出的成员函数声明为虚拟的，这样既可以避开 Visual C++ 潜在的问题，又可以让用户在必要时覆盖成员函数。

```
class CMyExportedClass
{
public:
    CMyExportedClass(void) : mdwValue(0) { }

    DLLFUNCTION virtual void setValue(long dwValue) { mdwValue = dwValue; }
    DLLFUNCTION virtual long getValue(void) { return mdwValue; }
    long clearValue(void);

private:
    long mdwValue;
};
```

将导出的成员函数声明为虚拟的以后，在应用程序中继承导出的类时，就像该类是在应用程序中声明的一样。

```
Class CmyApplicationClass : public CMyExportedClass
{
public:
    CmyApplicationClass (void) { }

    virtual void setValue (long dwValue);
    virtual long getValue (void);
};
```

#### 1.4.5 总结

从 DLL 中导出类是一种简单、功能强大的共享功能而又无需共享源代码的方式。它给应用程序提供了结构化 C++ 类的所有好处，应用程序可以使用、继承、重载这些类，同时让类的创建者可以安全地隐藏内部的函数和变量。



## 1.5 避免 DLL 困境

Herb Marselas, Ensemble Studios

hmarselas@ensemblestudios.com

**动**态链接库 (DLL) 是 Microsoft Windows 的一项强大特性, 有很多用途, 其中包括共享可执行代码以及抽象化设备差异。不幸的是, 由于其独立的特征, 依赖于 DLL 可能导致问题。如果一个应用程序依赖于一个用户计算机中没有的 DLL, 并试图运行该 DLL, 系统将显示消息“DLL Not Found”, 而这种消息对于普通用户来说是没有任何帮助的。如果应用程序启动时自动装载一个 DLL, 而用户的计算机中没有该 DLL, 则将无法指出该 DLL 是否有效 (至少对应用程序而言是这样的)。

用户安装和卸载程序时, 很容易使系统的 DLL 版本不正确。另外, 不同的 Windows 平台和服务软件包中的系统 DLL 也可能不同。在这些情况下, 用户可能收到晦涩的消息“DynaLink Error!” (DLL 中没有被链接的函数时), 甚至应用程序将崩溃。所有这些关于查找正确的 DLL 的问题常被称为“DLL 困境”。幸运的是, 避免陷入这种困境的方式有多种。

### 1.5.1 显式链接还是隐式链接

要避免 DLL 困境, 首先应确保用户的计算机中包含所需的 DLL, 且版本是应用程序能够使用的。在试图使用 DLL 的任何功能之前, 必须这样做。

通常, 我们在链接命令行 (link line) 指定 DLL 对应的 lib 文件, 来将 DLL 链接到应用程序。这被称为隐式 DLL 装载或隐式链接。通过链接 lib 文件, 当程序运行时, 操作系统将自动搜索并装载匹配的 DLL。这种方式假定 DLL 是存在的, Windows 将找到它, 且其版本正是程序要使用的。

Microsoft Visual C++ 也支持其他三种隐式链接方式。第一种是将 DLL 的 lib 文件直接包含到工程中, 这同将其加入到链接命令行中相同。第二种是, 如果工程包含一个构建 DLL 的子工程, 则默认情况下, 该 DLL 的 lib 文件将被自动链接到工程。最后一种是, 使用编译指令 #program comment (lib “libname”) 将 lib 链接到应用程序。

对于隐式链接和装载存在的问题, 解决办法是显式地装载 DLL。方法是, 不在链接命令行中进行到 DLL 的 lib 文件的链接, 并将所有将链接到库的编译指令 #program comment 删除。如果 Visual C++ 中的子工程将构建一个 DLL, 则修改其链接属性页, 选中其中的“Doesn't produce .LIB”选项。通过显式地装载 DLL, 代码将能够处理所有可能发生的错误, 确保

DLL 存在，确保所需的函数存在等等。

### 1.5.2 LoadLibrary 和 GetProcAddress

使用 lib 文件隐式地装载 DLL 后，可以在应用程序代码中直接调用函数，OS 装载程序将执行装载 DLL 和解析函数引用的所有工作。采用显式链接后，必须间接地通过一个手工解析的函数指针来调用函数。为此，必须使用函数 LoadLibrary 显式地装载包含该函数的 DLL，然后使用 GetProcAddress 来获得指向函数的指针。

```
HMODULE LoadLibrary(LPCTSTR lpFileName);
FARPROC GetProcAddress( HMODULE hModule, LPCSTR lpProcName);
BOOL FreeLibrary(HMODULE hModule);
```

LoadLibrary 将搜索指定的 DLL，如果找到，则将其装载到应用程序的进程空间中，并返回一个指向该新模块的句柄。然后，便可以使用 GetProcAddress 来创建指向游戏要使用的、该 DLL 中的每个函数的函数指针。当不再需要显式地装载的 DLL 时，应使用 FreeLibrary 来释放它。调用 FreeLibrary 后，相应的模块句柄将不再有效。

每个 LoadLibrary 调用都必须有一个对应的 FreeLibrary 调用，这是因为，无论是可执行文件或其他 DLL 隐式地装载 DLL，还是通过调用 LoadLibrary 显式地装载，Windows 的进程都会将该 DLL 的引用计数加 1。当调用 FreeLibrary 或装载该 DLL 的可执行文件或 DLL 被卸载时，引用计数将减 1。当 DLL 的引用计数为零后，Windows 将知道，此时可以安全地卸载该 DLL。

### 1.5.3 提防 DirectX

我们经常遇到的问题之一是，所需的 DirectX 组件版本没有安装或已损坏。为防止游戏出现这样的问题，我们显式地装载所需的 DirectX 组件。如果要隐式地链接 DirectX 8 中的 DirectInput，则必须在链接命令行中添加 dinput8.lib，并使用下面的代码：

```
IDirectInput8 *pDInput;
HRESULT hr = DirectInput8Create(hInstance, DIRECTINPUT_VERSION,
                                IID_IDirectInput8,
                                (LPVOID*) & pDInput, 0);

if (FAILED(hr))
{
    // handle error - initialization error
}
```

显式装载 DLL 时，代码多了两行，但现在应用程序能够防范 dinput8.dll 不存在或被损坏的情况。

```
typedef HRESULT (WINAPI* DirectInput8Create_PROC)
(HINSTANCE hinst, DWORD dwVersion, REFIID riidIthf, LPVOID* ppvOut,
```

```
LPUNKNOWN punkOuter);

HMODULE hDInputLib = LoadLibrary("dinput8.dll");
if (!hDInputLib)
{
    // handle error - DInput 8 not found. Is it installed incorrectly or at all?
}

DirectInput8Create_PROC diCreate;
diCreate = (DirectInput8Create_PROC)
GetProcAddress(hDInputLib, "DirectInput8Create");

if (!diCreate)
{
    // handle error - DInput 8 exists, but the function can't be found.
}

HRESULT hr = (diCreate)(hInstance, DIRECTINPUT_VERSION, IID_IDirectInput8,
                        (LPVOID*) &mDirectInput, NULL);

if (FAILED(hr))
{
    // handle error - initialization error
}
```

首先创建了一个反映函数 `DirectInput8Create` 的函数指针 `typedef`。然后使用 `LoadLibrary` 装载 DLL。如果成功地装载了 `dinput8.dll`，则使用 `GetProcAddress` 来找到函数 `DirectInput8Create`。如果找到，`GetProcAddress` 将返回一个指向该函数的指针；否则返回 `NULL`。然后对指针进行检查，以确保其有效。最后，通过该函数指针调用 `DirectInput8Create`，以初始化 `DirectInput`。

如果还需要使用该 DLL 中的其他函数，则需要为每个函数声明一个 `typedef` 函数指针和变量。也许，只需在使用 `GetProcAddress` 查找第一个函数时，检查函数指针是否为 `NULL` 即可。然而，添加更多的错误处理通常并不是什么坏事。检查每个 `GetProcAddress` 调用是否返回非 `NULL` 值是不错的选择。

#### 1.5.4 使用操作系统特有的特性

显式装载 DLL 可以解决的另一个问题是，应用程序要使用一个特定的 API 函数（如果该函数可用的话）。有许多以“Ex”结尾的扩展函数，这些函数在 Windows NT 或 Windows 2000 中是支持的，但在 Windows 95 或 Windows 98 中不可用。这些函数通常比原来的函数提供的信息或功能更多。

一个这样的例子是 `CopyFileEx`，它能够取消对长文件的复制操作。我们使用 `LoadLibrary` 装载 `kernel32.dll`，然后使用 `GetProcAddress` 查找该函数，而不是直接调用它。如果成功地装载了 `kernel32.dll`，并找到 `CopyFileEx`，则使用它；如果没有找到，则使用常规的 `CopyFile` 函数。这里必须避免的另一个问题是，`CopyFileEx` 实际上是头文件 `winbase.h` 中的一个 `#define`

替换, 针对 ASCII 或宽 Unicode 字符进行编译时, 将分别被替换为 CopyFileExA 和 CopyFileExW。

```
typedef BOOL (WINAPI *CopyFileEx_PROC) (LPCTSTR lpExistingFileName,
    LPCTSTR lpNewFileName, LPPROGRESS_ROUTINE lpProgressRoutine, LPVOID
    lpData, LPBOOL pbCancel, DWORD dwCopyFlags);

HMODULE hKernel32 = LoadLibrary("kernel32.dll");
if (!hKernel32)
{
    // handle error - kernel32.dll not found. Wow! That's really bad
}

CopyFileEx_PROC pfnCopyFileEx;
pfnCopyFileEx = (CopyFileEx_PROC) GetProcAddress(hKernel32, "CopyFileExA");

BOOL bReturn;
if (pfnCopyFileEx)
{
    // use CopyFileEx to copy the file
    bReturn = (pfnCopyFileEx)(pExistingFile, pDestinationFile, ...);
}
else
{
    // use the regular CopyFile function
    bReturn = CopyFile(pExistingFile, pDestinationFile, FALSE);
}
```

LoadLibrary 和 GetProcAddress 也可用于游戏 DLL。一个这样的例子是 Ensemble 工作室正在开发的游戏引擎中的图形支持。在该游戏引擎中, 对 Direct3D 和 OpenGL 的图形支持被放在不同的 DLL 中, 必要时可以显式地装载这些 DLL。如果需要 Direct3D 图形支持, 则使用 LoadLibrary 装载支持 Direct3D 的 DLL, 并使用 GetProcAddress 找到导出的函数。这种设计使得可执行主文件无需隐式地链接 d3d8.lib 或 opengl32.lib。

然而, 支持 Direct3D 的 DLL 隐式地链接 d3d8.lib, 而支持 OpenGL 的 DLL 隐式地链接 opengl32.lib。这种由有可执行主文件显式地装载游戏自己的 DLL, 每个图形子系统隐式地进行装载的方式解决了几个问题。首先, 如果装载库失败, 则很可能是特定的图形子系统文件不存在或已损坏, 主程序将能够妥善地处理错误。这种方式解决的另一个问题 (这种问题通常出现在 OpenGL, 而不是 Direct3D 中) 是, 如果引擎要显式地链接到 OpenGL, 则对于它使用的每个 OpenGL 函数, 都需要一个 typedef 和函数指针。隐式地链接到支持的 DLL 解决了这种问题。

### 1.5.5 总结

显式链接能够防范许多常见的 DLL 问题, 这些问题通常出现在 Windows 中, 其中包含 DLL 不存在或 DLL 版本与应用程序不兼容。虽然它不是万能药, 但至少能够对应用程序进

行控制，从而妥善地对所有错误进行处理，而不是显示一个晦涩的错误消息或导致应用程序立刻崩溃。





## 1.6 动态类型信息

Scott Wakeling, Virgin Interactive Ltd.

scott@chronicreality.com

随着开发人员采用面向对象技术,给游戏提供动力的系统越来越灵活,也更为复杂。现在,这种系统经常包含很多不同的类型和类——空前地超过 1000 个。在游戏中处理这么多不同的类型本身就是一种挑战。类型可以是类、结构或标准的数据类型。本文讨论如何查询类型同其他类型的关系——出于查询或调试的目的,在运行阶段访问有关类型的信息,从而高效地管理类型。文章的最后提供了一种支持永久性对象的方法以及如何扩展该方法。

### 1.6.1 动态类型信息类简介

为高效地利用类型的功能,我们求助于一种特殊的类:动态类型信息(DTI)类。这个类将存储有关对象或结构的各种信息,其最简单的实现如下:

```
class dtiClass
{
private:
    char* szName;
    dtiClass* pdtiParent;
public:
    dtiClass();
    dtiClass( char* szSetName, dtiClass* pSetParent );
    virtual ~dtiClass();

    const char* GetName();
    bool SetName( char* szSetName );

    dtiClass* GetParent();
    bool SetParent( dtiClass* pSetParent );
};
```



要在引擎中安装 DTI,所有的类都应将其一个静态成员。这个类让我们在调试时能够访问类的名称以及查询父类的 dtiClass 成员。在继承树中,所有的类都应包含这样的成员,这样所有的游戏对象都能够访问有关自身及其父类的信息。本书附带光盘中提供了 dtiClass 类的实现代码。

## 1.6.2 暴露和查询 DTI

接下来按前面的介绍实现一个非常简单的类树，以使用 DTI。下面的代码片段包含一个帮助定义静态 `dtiClass` 成员的宏、一个基本的根类，并对类的类型信息进行了初始化。

```
#define EXPOSE_TYPE \
    public: \
        static dtiClass Type;

class CRootClass
{
public:
    EXPOSE_TYPE;
    CRootClass() {};
    virtual ~CRootClass() {};
};

dtiClass CRootClass::Type( "CRootClass", NULL );
```

通过在所有的类定义中包含 `EXPOSE_TYPE` 宏，并正确地初始化静态成员 `Type`，我们向在游戏引擎中安装动态类型信息迈出了第一步。我们将类名和一个指针传递给父类的 `dtiClass` 成员，其余的工作将由 `dtiClass` 的构造函数来完成——相应地设置成员 `szName` 和 `pdtiParent`。

现在，便可以在运行阶段，为调试与类型相关的问题——如保存或装载游戏，而查询对象的类名。后面将对此做更详细的介绍，下面为获取类名的代码：

```
// Let's see what kind of object this pointer is pointing to
const char* szGetName = pSomePtr->Type.GetName();
```

前面的范例将 `NULL` 作为父类名传递给 `dtiClass` 的构造函数，因为这是根类。对于从其他类派生而来的类，只需指定父类的名称即可。例如，要从根类派生出子类，则其定义如下：

```
class CChildClass : public CRootClass
{
    EXPOSE_TYPE;
    // Constructor and virtual Destructor go here
};

dtiClass CChildClass::Type( "CChildClass", &CRootClass::Type );
```

现在有一个类树，我们不但能够访问类的名称，还能访问其父类的名称——条件是使用了 `EXPOSE_TYPE` 宏暴露其类型。获得父类名称的代码如下：

```
// Let's see what kind of class this object is derived from
char* szParentName = pSomePtr->Type.GetParent()->GetName();
```

有了包含 DTI 的类树，并知道如何在运行阶段使用这些信息来查询类及其父类的名称后，

便可以实现一个方法，用于保障强制类型转换，或查询对象的祖先（通用类型）。

### 1.6.3 继承的含义是“是一个”

面向对象技术提供了继承的能力。伴随继承而来的是多态——在任何时候，所有的对象都可被视为某种类型。在游戏编程中，很多情况下都通过多态来安全、高效、动态地处理多种类型的对象。这意味着对对象进行强制类型转换之前，必须确保其类型是兼容的，从而防止未定义的行为发生。另外，还需要能够在运行阶段（而不是编译阶段）快速、便捷地检查对象与哪种类型兼容。

假设游戏涉及到大量不同的机器人，其中的一些纯粹是电子化的，而有一些包含机械部件，并通过燃料来驱动。现在假设玩家拥有一种武器，它对付纯粹的电子化机器人非常有效，但对付机械化机器人则不那么有效。定义这些机器人的类很可能是同一种基本类型，即它们是从同一个一般机器人基类派生而来的，但覆盖了某些函数，并新增了一些属性。为处理各种不同的子类，需要查询它们的父类。我们可以对前面介绍的 `dtiClass` 进行扩展，以提供这样的例程。我们将这个新的成员函数命名为 `IsA`，因为继承的意思是“是一个”。该函数的代码如下：

```
bool dtiClass::IsA( dtiClass* pType )
{
    dtiClass* pStartType = this;
    while( pStartType )
    {
        if ( pStartType == pType )
        {
            return true;
        }
        else
        {
            pStartType = pStartType->GetParent();
        }
    }
    return false;
}
```

如果需要知道某个机器人子类是否是从特定的根类派生而来的，只需调用该对象的 `dtiClass` 成员的 `IsA` 函数，并将根类的静态成员 `dtiClass` 传递给该函数即可。下面是一个简短的范例：

```
CRootClass* pRoot;
CChildClass* pChild = new CChildClass();
if ( pChild->Type.IsA( &CRootClass::Type ) )
{
    pRoot = (CRootClass*)pChild;
}
```



我们从中可以知道，IsA 指出当前对象是否是从给定的基类直接或间接地派生而来的。当然，我们可以根据这种信息执行安全的强制类型转换操作，就像上面的范例那样。另外，也可以通过这种检查来过滤掉给定区域中的某种游戏对象，因为它们经不起某种武器或效果的打击。如果经常需要执行安全的强制类型转换，则可以将下面的函数加入到根对象中，以简化工作。这里只提供了它的定义和一个简短的例子，有关该函数的实现请参阅附带光盘。

```
// SafeCast member function definition added to CRootClass
void* SafeCast( dtiClass* pCastToType );

// How to simplify the above operation
pRoot = (CRootClass*)pChild->SafeCast( &CRootClass::Type );
```

如果强制类型转换是不安全的（即类型是不相关的），则值为空，而 pRoot 将为 NULL。

#### 1.6.4 处理通用对象

回到前面的游戏范例，让我们考虑如何高效地处理多种不同的机器人。解决方案很简单：可以利用多态，将指向它们的指针存储在一个大型的通用基类指针数组中。即使是更具体的机器人，如 CRobotMech（它是从 CRobot 派生而来的），其指针也可以存储在这种数组中，因为多态使得对于任何类型要求，派生类型总能满足。现在有一个大型的游戏对象数组，这些对象都被存储为指向给定基类的指针。因此可以安全地遍历它们——这可能是调用各个对象的虚函数或更具体的函数（覆盖函数）。这样，我们便迈向了以快速、安全、通用的方式处理大量的游戏对象之路。

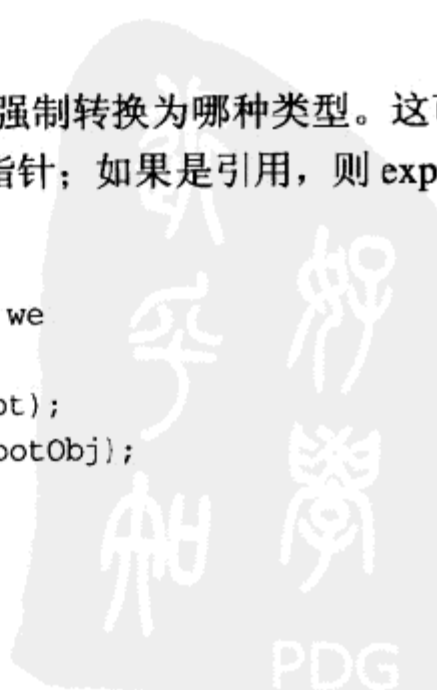
在运行阶段类型信息解决方案中，使用了函数 IsA 和 SafeCast，它们分别用于查询对象的通用类型和安全地沿类树向上进行强制类型转换。这常被称为向上强制类型转换，它使我们迈向了以快速、安全、通用的方式处理大量的游戏对象之路。问题的另一部分是向下强制类型转换——将指向通用基类的指针安全地转换为指向具体子类的指针。要遍历根类指针数组，并检查其中的每个指针是否指向特定的子类，需要利用 C++ 引进的动态类型转换运算符。

动态类型转换运算符用于在多态类型之间进行转换，它既安全又富含信息。它甚至能够返回有关试图进行转换的信息。该运算符的格式如下：

```
dynamic_cast< type-id >(expression)
```

其中第一个参数是必须传递的，它指出要将 expression 强制转换为哪种类型。这可以是指向类的指针或引用。如果是指针，则 expression 也必须是指针；如果是引用，则 expression 必须是一个可修改的左值 (l-value)。下面是两个例子：

```
// Given a root object (RootObj), or pointer (pRoot) we
// can down-cast like this
CChildClass* pChild = dynamic_cast<CChildClass*>(pRoot);
CChildClass& ChildObj = dynamic_cast<CChildClass&>(RootObj);
```





要使用这些扩展的类型转换运算符，必须在编译器设置中启用嵌入的运行阶段类型信息（在 Microsoft Visual C++ 中，使用 /GR 开关）。如果请求的强制类型转换无法进行（例如根类指针并非指向某个派生类），运算将失败，结果为 NULL。因此，在上述代码片段中，如果 pRoot 指向的是一个 CRootClass 对象，则 pChild 的值将为 NULL。如果转换 ChildObj 失败，将引发异常，因此该转换应放在一个 try/catch 语句块中（范例请参见附带光盘）。



运算符 dynamic\_cast 让我们能够确定指针的真正类型。假设我们要遍历特定半径范围内的所有机器人，并确定哪些机器人是机械化的，能够经受住某种武器的打击。对于一个给定的 CRobot 指针数组，可以遍历其中的每个指针，对其执行动态类型转换，并检查哪些转换是成功的，哪些转换结果为 NULL，从而挑选出真正的机械化机器人。现在，我们也能够安全地进行向下类型转换，完成了运行阶段类型信息解决方案。附带光盘中包含一个更复杂的使用动态类型转换运算符的范例。

### 1.6.5 实现永久性类型信息



可以在运行阶段确定对象的身份，并高效地管理它们后，我们便可以实现一个永久性对象解决方案，从而增加与类型相关的功能，并轻松地处理诸如游戏保存或对象仓库等问题。首先，需要一个二进制存储的实现框架，用于保存对象数据。附带光盘中包含一个名为 CdtiBin 的实现范例。

成员函数很多，其中比较重要的两个是 Stream 和友元运算符 <<，它们用于装载语言的基本数据类型或将其写入到磁盘中。对于要永久化的每种基本类型，都需要添加一个运算符。当 Stream 被调用时，它将根据 m\_bLoading 和 m\_bSaving 的值，从文件中读取数据或将数据写入到文件中。

为让类知道如何使用对象仓库，需要添加 Serialize 函数，如下所示：

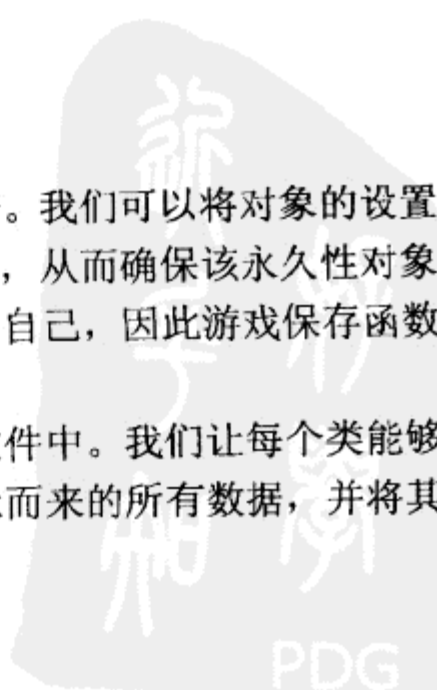
```
virtual void Serialize( CdtiBin& ObjStore );
```

这是一个虚函数，对于有额外数据的子类，都需要覆盖这个函数。例如，如果在 CRootClass 的基础上新增一个 int 成员，则 Serialize 函数将如下：

```
void CRootClass::Serialize( CdtiBin& ObjStore )
{
    ObjStore << iMemberInt;
}
```

另外，我们还需要为 CdtiBin 和 int 成员提供友元运算符。我们可以将对象的设置写入到文件中，以后可以装载这些设置，并使用它们来设置新对象，从而确保该永久性对象解决方案可用于游戏保存函数中。各种类型的对象都知道如何保存自己，因此游戏保存函数实现起来将容易得多。

然而，子类需要将其自身的数据及父类的数据写入到文件中。我们让每个类能够访问其父类的 Serialize 例程，从而避免需要程序员查看从父类继承而来的所有数据，并将其加入到





每个类的 `Serialize` 成员函数中。这样，子类可在写入（装载）自身的数据之前，写入（装载）继承的数据。为此，我们使用了 `DECLARE_SUPER` 宏：

```
#define DECLARE_SUPER(SuperClass) \
    public: \
        typedef SuperClass Super;

class CChildClass
{
    DECLARE_SUPER(CRootClass);
    // ...
}
```

通过让类可以调用其直接父类的函数，我们进一步扩展了类型解决方案，提高了类树的扩展性。

`CRootClass` 无需声明其父类，因为它没有父类，因此其成员函数 `Serialize` 只需处理自身的数据。下面的代码演示了 `Class::Serialize` 函数在处理自己的数据之前，调用 `CRootClass::Serialize`。

```
void CChildClass::Serialize( CdtiBin& ObjStore )
{
    Super::Serialize( ObjStore );

    ObjStore << fMemberFloat << iAnotherInt;
}
```

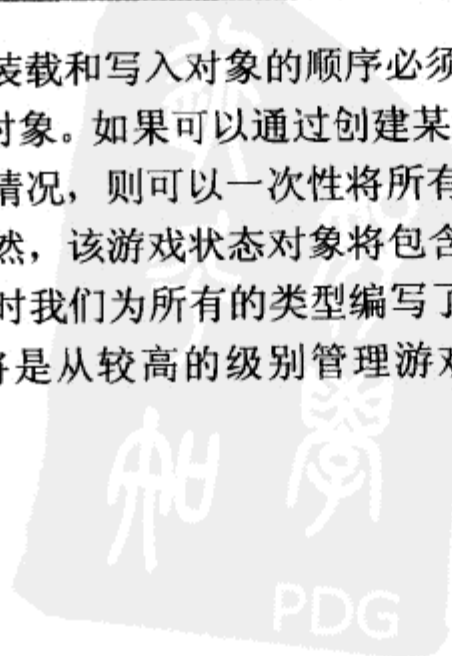


为支持上述代码，我们为数据类型 `float` 添加了一个友元运算符。注意，属性的保存和装载顺序总是相同的。附带光盘中的代码演示了如何创建二进制存储，将对象的数据写入到磁盘，然后使用这些数据来设置对象的属性。

只要在保存和装载时序列化对象类型的顺序相同，便可永久地保持对象的属性不变。可在 `CdtiBin` 类中新增友元运算符，以支持基本类型。如果要用户定义的结构作为类成员，只需编写一个处理该结构的运算符即可。这样，引擎中的所有对象和类型都知道如何存储和读取自己。

### 1.6.6 将永久性类型信息用于游戏保存数据库中

正如前面指出的，装载和写入对象的顺序必须相同。最快捷、最简单的方法是，只保存一个对象，然后装载该对象。如果可以通过创建某种游戏状态对象（它知道如何序列化自身）来定义游戏任何时候的情况，则可以一次性将所有的游戏数据写入到文件中，并在以后任何时候读取这些数据。显然，该游戏状态对象将包含一个对象数组。只要该定制的数组类型知道如何序列化自身，同时我们为所有的类型编写了正确的 `CdtiBin` 运算符，将万事大吉。这样，保存和装载游戏将从较高的级别管理游戏——创建容器类，并在必要时调用一个 `Serialize` 例程。





### 1.6.7 结论

---

上述解决方案还有很多可改进的地方。支持多重继承将不再困难，可以在静态成员 `dtiClass` 中存储一个包含所有父类的数组（而不是一个父类指针），在一个合适的宏中（或通过扩展 `dtiClass` 类的构造函数）指定父类的数目和父类的类型。对象标记（flagging）系统也将很有用，它让我们能够处理特殊情况，如抽象基类或包含在其他类中的类（被包含的类）。

### 1.6.8 参考文献

---

[Meyers98] Meyers, Scott D., *Effective C++ 2<sup>nd</sup> Edition*, Addison-Wesley, 1998.

[Wilkie94] Wilkie, George., *Object-Oriented Software Engineering*, Addison-Wesley, 1994.

[Eberly00] Eberly, David H., *3D Game Engine Design*, Morgan Kaufman, 1999~2000.

[Wakeling01] Wakeling, Scott J., "Coping with Class Trees," available online at [www.chronicreality.com/articles](http://www.chronicreality.com/articles), March 12, 2001.



## 1.7 用于通用 C++ 成员访问的属性类

Charles Cafrelli  
skywise@iquest.net

几乎每个游戏都有一组独特的对象，而在每个项目中，操纵这些对象的代码必须从空白开始编写。例如，游戏编辑器，其用途很简单：创建、放置、显示和编辑对象的属性。对象的创建总是因游戏而异，这也可以由类工厂（factory）来处理。对象的放置随可视化引擎而异，这使得重用很困难，即使假定可以将对象放置到地图上。在有些情况下，可在不同的游戏中重用可开关的通用地图编辑器。因此，从理论上说，可以开发一个可重用的核心编辑器模块，而无需针对不同的项目重新编写代码。然而，由于每个游戏都有其独特的对象，在不重写编辑器代码的情况下，编辑器如何知道应显示哪些内容供用户进行编辑呢？

我们需要一个通用的对象接口，通过它可以访问类中的数据。Borland C++ Builder 提供了一个非常棒的 C++ 声明类型——property，能够完成这样的工作，然而，这是一种专用扩展，在 Borland C++ 外不可用。有趣的是，由 Borland C++ Builder 开发人员创建的 Microsoft 的新的编程语言 C# 也提供了这样的特性。Microsoft 的 COM 接口允许在运行阶段查询对象的成员，但它要求将对象与 COM 接口绑定起来，这使得其可移植性不如 C++ 高。因此，您只能创建自己的解决方案，其量级比 COM 轻，而可移植性比专用的 C++ 扩展高。这样，对于诸如游戏编辑器等代码模块，只需编写一次，并可在多个引擎中使用。

### 1.7.1 代码



该接口被划分为两个类：Property 类和 PropertySet 类。前者是一个存储数据的容器，其中包含一个共用体、一个枚举和一个字符串。它们分别是指向不同数据类型的指针、数据类型和属性名。有关该接口的完整源代码，请参阅附带光盘。

```
class Property
{
protected:
    union Data
    {
        int* m_int;
        float* m_float;
```

```

        std::string* m_string;
        bool* m_bool;
    };

    enum Type
    {
        INT,
        FLOAT,
        STRING,
        BOOL,
        EMPTY
    };

    Data m_data;
    Type m_type;
    std::string m_name;

protected:
    void EraseType();
    void Register(int* value);
    void Register(float* value);
    void Register(std::string* new_string);
    void Register(bool* value);

public:
    Property();
    Property(std::string const& name);
    Property(std::string const& name, int* value);
    Property(std::string const& name, float* value);
    Property(std::string const& name, std::string* value);
    Property(std::string const& name, bool* value);
    ~Property();

    bool SetUnknownValue(std::string const& value);
    bool Set(int value);
    bool Set(float value);
    bool Set(std::string const& value);
    bool Set(bool value);

    void SetName(std::string const& name);
    std::string GetName() const;

    int GetInt();
    float GetFloat();
    std::string GetString();
    bool GetBool();
};

```

上述代码范例使用和存储的是基本数据类型，但是可以对上述代码进行扩展，以处理任何数据类型。属性存储的是一个指针，它指向原来的数据。属性并没有声明对象，也没有分配内存，因此操纵属性时，实际上处理的是原来的数据。通过 Set 函数设置值时，将自动定义属性的类型。

属性是通过 PropertySet 类来创建和操纵的。PropertySet 类包含一系列的注册属性、注册方法和查找方法。

```
class PropertySet
{
protected:
    HashTable<Property> m_properties;

public:
    PropertySet();
    virtual ~PropertySet();

    void Register(std::string const& name, int* value);
    void Register(std::string const& name, float* value);
    void Register(std::string const& name, std::string* value);
    void Register(std::string const& name, bool* value);

    // look up a property
    Property* Lookup(std::string const& name);

    // get a list of available properties
    bool SetValue(std::string const& name, std::string* value);
    bool Set(std::string const& name, std::string const& value);
    bool Set(std::string const& name, int value);
    bool Set(std::string const& name, float value);
    bool Set(std::string const& name, bool value);
    bool Set(std::string const& name, char* value);
};
```



PropertySet 是围绕一个 HashTable 对象组织的，该对象使用简单的哈希表算法来组织所有被存储的属性。HashTable 对象本身是一个模板，可用来处理不同的对象，其代码请参见附带光盘。

我们以 PropertySet 类为基类，派生出游戏对象：

```
class GameObject : public PropertySet
{
    int m_test;
};
```

对于任何需要以公有方式暴露（或被其他对象使用）的属性，都需要进行注册，这通常是在构造时完成的。例如：

```
Register("test_value",&m_test);
```

调用对象可以使用 Lookup 方法来访问注册的数据：

```
void Update(PropertySet& property_set)
{
    Property* test_value_property=
        property_set.Lookup("test_value");
```



```
int test_value = test_value_property->GetInt();  
// blah  
}
```

由于现在所有对象的类型都为 `PropertySet`，且通常被存储在主更新列表中，因此将列表指针放到游戏编辑器中进行处理将很简单。新的派生对象类型只需注册其新增的属性，以便编辑器能够进行处理；而无需编写其他的代码，因为编辑器并不关心派生类型。有时候，在属性名（如 `Type`）中指定类型，以帮助用户可视化地编辑对象将很有用。另外，使属性为必须的，以便编辑器能够以树的方式显示属性列表也很有用。

这也提供了另一个好处，即降低了数据同其名称之间的关联性。例如，在内部，数据可以名为 `m_colour`，但暴露为“color”。

### 1.7.2 其他用途

这些类是按照同心环设计理论进行设计的。如果没有 `Property` 类，`PropertySet` 类将无法使用。然而，`Property` 类可以单独使用，也可以与其他设置类型（如 `MultiMatrixedPropertySet`）一起使用，而无需重新编写；`PropertySet` 类中的 `HashTable` 也是如此。与包含用于处理各种可能用途的方法的大型类相比，用于特定用途的小型类的可重用性更高。

`Property` 类也可用于暴露方法，这样就可以从外面通过函数指针调用这些方法。只需再编写少量的代码，它还可以用作保存状态，从而实现游戏保存特性。另外，它还可用于通过网络在对象之间发送消息。添加 `Send(std::string xml)` 和 `Receive(std::string xml)` 方法后，`Property` 还可以对包含属性值（或需要修改的属性值）的 XML 消息进行编码和解码。也可以将 `Property/PropertySet` 类重写为模板，以支持不同的属性类型。

使用 `get` 和 `set` 方法来隔离数据，可以在内部存储格式和外部格式之间进行转换。这样，将无需知道属性的数据类型，使代码更为通用，但代价是速度将稍慢些。

### 1.7.3 推荐读物

Fowler, Martin, Kent Beck, John Brant, William Opdyke, Don Roberts, *Refactoring*, Addison-Wesley, ISBN: 0201485672.

Gamma, Erich, Richard Helm, Ralph Johnson, John Vlissides, Grady Booch, *Design Patterns*, Addison-Wesley, ISBN: 0201633612.

Lakos, John, *Large-Scale C++ Software Design*, Addison-Wesley, ISBN: 0201633620.

McConnell, Steve C., *Code Complete: A Practical Handbook of Software Construction*, Microsoft Press, ISBN: 1556154844 (anything by McConnell is good).

Meyers, Scott, *Effective C++: 50 Specific Ways to Improve Your Programs and Design (2nd Edition)*, Addison-Wesley, ISBN: 0201924889.

Meyers, Scott, *More Effective C++: 35 New Ways to Improve Your Programs and Designs*, Addison-Wesley, ISBN: 020163371X.

## 1.8 一个游戏实体工厂

François Dominic Laramée

francoislaramée@videotron.ca

近年来,事实证明,对于游戏开发而言,脚本语言的价值是无法估量的。通过将改进游戏实体行为的代码从代码核心中分离出来,脚本语言使关卡设计器(level designer)能够脱离代码的编译—执行循环,将游戏测试和调整的速度提高了几个数量级,并让高级程序员能够将时间用于完成更复杂的任务。

然而,为使数据驱动的开发模式得以顺利地进行,游戏引擎必须提供灵活的实体构造和组装服务,这样脚本语言才能为各个实体提供不同的作战策略、反应行为和其他参数。本文旨在描述一个 C++ 类层次和一套技术,以支持引擎方面的数据驱动开发模式:

- 将逻辑行为和音频视频行为分开。必要时,一个 Door 类能够支持这一概念的很多变体,而无需关心动画序列的大小和关键帧数目。
- 快速开发。定义好基本行为库后(需要的时间很少),便可以在框架中添加新的游戏实体,这只需很少的代码,所需时间通常不超过 15 分钟。
- 避免代码重复。通过在运行阶段将行为组装到新实体中,该框架避免了使用脚本语言时,代码过长(相对于 C/C++)的问题。

本文以模式的方式描述了几种技术,有关模式的详细信息,请参阅被称作“四人帮(Gang of Four's)”编写的 *Design Patterns*[GoF94]一书。

### 1.8.1 组件

本文的讨论围绕着 3 个主要的组件: flyweight 对象、行为类和一个对象工厂方法。我们将依次介绍这些组件,然后讨论它们如何协同工作,使引擎能够提供数据驱动开发模式所需的服务。最后,将讨论如何使该系统更灵活(代价是代码更复杂)——如果项目需要一种功能齐备的脚本语言的话。

### 1.8.2 flyweight 类、行为类和导出类

首先,我们来介绍一下在本框架中,游戏实体包含的 3 种类之间的区别: flyweight 类、行为类和导出类

- flyweight 类是实体的外观。在代码中,实体同其 flyweight 类之间



的关系是通过对象组合实现的：实体包含一个指向 flyweight 类的指针，它使用该 flyweight 类来表示其外观。

- 行为类定义了对象如何同游戏世界中的其他对象进行交互。行为类是通过传统的继承层次实现的，Entity 类被用作其他所有类的抽象超类。

- 导出 (exported) 类指定了对象如何向世界描绘自己。这个类不是必须的，使用它旨在提供方便。导出类被实现为一个枚举常量，让实体在其生命周期内能够将自己宣告为多个不同的对象类。

下面依次介绍这些类。

### 1.8.3 flyweight 对象

[GoF94]将 flyweight 对象定义为剔除了场境的对象，因此可在不同的情形下使用；换句话说，它是其他对象的模板或模型。就游戏实体而言，flyweight 友好的信息包括：

- 媒体内容：音响效果、3D 模型、纹理、动画文件等；
- 控制结构：有限状态机定义、脚本等。

正如您看到的，这可以是除实体当前状态（位置、生命值、FSM 状态）之外的任何信息。因此，在游戏中，flyweight 并不适合，因为 flyweight 可能占用数兆的内存，而场境信息却小得甚至能容纳进一个烤面包机的核心内存。

### 1.8.4 SAMMy，您在哪里？

游戏实体的有限状态机主要处理的是动画循环，决定什么时候播放声音等。例如，街机中当玩家被杀时，可能进入复活 (recurrecting) 状态，并被标记为不会受到伤害，同时播放复活动画；否则，妖怪将在附近逗留，并在动画的每一帧中再下杀手，直到玩家获得控制权为止。我将 flyweight 对象中处理这种情况的部分称作状态和媒体管理器 (State And Media Manager)，简称为 SAMMy：

```
class StateAndMediaManager
{
    // The various animation sequences available for the family
    // of entities
    AnimSequenceDescriptionStruct * sequences;
    int numAnimSequences;

    // A table of animation sequences to fire up when the entity's FSM
    // changes states out of SAMMy's control
    int * stateToAnimTransitions;
    int numStateToAnimTransitions;

public:
    // Construction and destruction
    // StateAndMediaManager is always constructed by its owner entity,
    // which is in charge of opening its description file. Therefore,
```

```

// the only parameter the constructor needs is a reference to an
// input stream from which to read a set of animation sequence
// descriptions.
StateAndMediaManager() : sequences( 0 ), numAnimSequences( 0 ),
    numStateToAnimTransitions( 0 ), stateToAnimTransitions( 0 ) {}
StateAndMediaManager( istream & is );
virtual ~StateAndMediaManager();
void Cleanup();

// Input-output functions
void Load( istream & is );
void Save( ostream & os );

// Look at an entity's current situation and update it according
// to the description of its animation sequences
void FC UpdateEntityState( EntityStateStruct * state );

// If the entity's FSM has just forced a change of state, the media
// manager must follow suit, interrupt its current animation
// sequence and choose a new one suitable to the new FSM state
void FC AlignWithNewFSMState( EntityStateStruct * state );
};

```



通常, SAMMy 是使用实体制作工具创建的, 需要时从文件装载到引擎中。出于简化的目的, 附加光盘中的范例被创建一个文本文件。

可以根据需要, 使 SAMMy 强大而通用。从理论上说, SAMMy 可以负责所有的控制功能: 启动脚本、修改策略等。然而, 如果这样做, 它将难以使用, 且需要耗费大量的心血。我们将在行为类层次中实现大部分的高级控制结构, 这样只需很少的代码便可实现。共同完成任务的副作用在于, 单个行为类 (如 SecurityGuard、Door 或 ExplosionFx) 将能够基于多个相关的 flyweight 来处理实体, 这样系统将更灵活。

### 1.8.5 行为类层次

这是实体所属的 C++ 类。层次结构 (至少) 有两层:

- 抽象基类 Entity, 它定义了接口和共性;
- 具体的子类, 它是从 Entity 派生而来的, 并实现了实际的对象。

Entity 的接口如下:

```

class Entity
{
    // Some application-specific data

    // Flyweight and Exported Class information
    int exportedClassID;
    StateAndMediaManager * sammy;

```



```

public:
    // Constructors
    ...

    // Accessors
    int GetExportedClass() { return exportedClassID; }
    StateAndMediaManager * GetFlyweight() { return sammy; }
    void SetExportedClass( int newval ) { exportedClassID = newval; }
    void SetFlyweight( StateAndMediaManager * ns ) { sammy = ns; }

    // Factory method
    static Entity * EntityFactory( int exportedClassRequested );

    virtual Entity * CloneEntity() = 0;
    virtual bool Updateself()
    {
        // Do generic stuff here; looping through animations, etc.
        return true;
    }

    virtual bool HandleInteractions( Entity * target ) = 0;
};

```

从上可知，在层次结构中添加一个新类需要做的工作很少：除构造函数外，最多只需覆盖基类的三个方法（也可能是两个）。

- **Clone()**只是调用复制构造函数；
- **UpdateSelf()**运行实体的内部机制。对于有些实体，这可能很简单，只是调用 SAMMY 中的相应方法来更新当前的动画帧；而对于其他实体，如玩家角色，这可能更复杂得多。
- **HandleInteractions()**在实体决定是否需要根据其他对象的行为和位置而修改其内部状态时被调用。其默认实现为空，换句话说，对象在弄虚作假。



附带的光盘中有 Entity 子类的范例，包含有一个玩家角色驱动器。

### 1.8.6 使用模板方法模式来完成行为任务

如果游戏中有多个行为类，而它们之间很容易区别，则一种名为模板方法模式[GoF4]的技术将很有帮助。这包括一个基类方法，它定义了一种算法，该算法在通过多态调用的子类方法中实现。

例如，所有的 PlayerEntity 对象查询输入设备，并在 UpdateSelf()方法中移动自己，但如何移动则取决于使用的输入设备、角色类型（FleetingRogue 比 OneLeggedBuddha 走得更快）等。因此 PlayerEntity 类将 UpdateSelf()定义为一个虚方法，该方法只在子类中实现。

```

class PlayerDevice : public Entity
{

```



```
// ...
void UpdateYourself();
void QueryInputDevice() = 0; // No implementation in PlayerDevice
};

class JoystickPlayerDevice : public PlayerDevice
{
    // ...
    void QueryInputDevice();
};

void PlayerDevice::UpdateYourself()
{
    // do stuff common to all types of player devices
    QueryInputDevice();

    // do more stuff
}

void JoystickPlayerDevice::QueryInputDevice()
{
    // Do the actual work
}
```

被正确使用时，模板方法模式将最大限度地减少编程中的剪切—粘贴操作，在软件工程中，这种“反模式”行为将导致灾难[Brown98]！

### 1.8.7 导出类

导出类是一种提供方便的技术，让您能够使实体的内部状态信息对脚本编写者透明。例如，假设要编写 Pac-Man 的 `HandleInteractions()` 方法。您可能首先检查它是否与某个幽灵发生碰撞，如果发生碰撞，则发生的情况取决于幽灵是否害怕（被吃）——被吃后返回基地（什么事也没发生）或猎杀（Pac-Man 死亡）：

```
void PacMan::HandleInteractions( Entity * target )
{
    if ( target->GetClass() == GHOST && target->GetState() == AFRAID )
    {
        score += 100;
        target->SendKillSignal();
    }
}
```

然而，如果需要给 Ghost 对象添加状态，情况将如何呢？例如，您可能希望幽灵的 SAMMY 包含一个“恐惧”的动画循环，每当 Pac-Man 跑过一个电线杆，该动画便运行 1 秒钟。如果您添加了 `GettingScared` 状态，SAMMY 将处理这种情况。然而，此时您需要在事件处理程序中检测 `GettingScared` 状态。

```

void PacMan::HandleInteractions( Entity * target )
{
    if ( target->GetClass() == GHOST &&
        ( target->GetState() == AFRAID ||
          target->GetState() == GETTINGSCARED ) )
    {
        ...
    }
}

```

这种方法很笨拙，当您在 SAMMy 中添加状态（从外面看来，这没有引入任何新的东西），需要相应地修改事件处理程序。我们引入导出类的概念，它是 Entity 对象中一个可查询的值，描述了该对象如何向世界宣称自己。这个值是在 UpdateSelf() 中维护的，它可以有多种形态。出于简化的目的，我们假设它是枚举列表中的整型常量之一：

```
enum { SCAREDGHOST, ACTIVEGHOST, DEADGHOST };
```

对于瞬间的、与动画相关的状态（如 GettingScared），无需导出任何信息。对于 Pac-Man 而言，幽灵可能处于死亡、活动或恐惧状态。两帧以前它开始恐惧、已经吓坏了还是逐渐平静下来是无关紧要的。通过使用导出类，而不是内部的 FSM 状态，Ghost 对象可以宣称自己是一个已死的、活动的或恐惧的幽灵。从外界看来，这相当于转换为 3 种不同的实体类，其代价是需要添加一个整数。这样，Pac-Man 的事件处理程序将如下：

```

void PacMan::HandleInteractions( Entity * target )
{
    if ( target->GetExportedClass() == SCAREDGHOST )
    {
        score += 100;
        target->SendKillSignal();
    }
}

```

结果是，代码更为清晰，维护起来也更容易，因为实体可能的导出类通常数目很小，且易于提早确定。而在开发过程中，随着新的外观和效果被加入到对象中，SAMMy 的 FSM 将急剧增大。

### 1.8.8 实体工厂

至此，我们拥有了所需的所有工具，可将它们用于创建对象。

关卡（level）文件包含多个实体声明，其中每个声明都指出了实体的行为类、flyweight 类、导出类、起始位置、速度以及类特定的参数（如妖怪的杀伤力、炸弹的启动时间、玩家工具箱的容量等）。为了简化，针对关卡设计器，我们做以下可行的假设：虽然行为类可能将自己宣称为多种导出类，但导出类只与一个行为类相关联。这样，在关卡文件中无需指定行为类，而且类层次与工具和关卡设计器分离开来了。下面是关卡文件中的代码片段：

```
<ENTITY Blinky>
```

```
<EXPORTEDCLASS ActiveGhost>
<XYZ_POSITION ...>
<PARAMETERS ...>
</ENTITY>
```

接下来在 **Entity** 类中添加一个工厂方法 (factory method)。

工厂是一个函数，其功能是根据需要构造任意数目类的对象实例。工厂将处理行为类层次中任何成员（具体）的请求。就编程而言，工厂方法很简单：

- 它有一个注册表，该注册表描述了已装载到游戏中的 flyweight 以及属于每个行为类的导出类。
- 它在必要时装载 flyweight。被请求创建不可见的 flyweight 类的实例时，首先创建并装载该 flyweight 的 SAMMy 对象。
- 被请求创建已装载的 flyweight 类的另一个实例时，工厂将克隆已有的对象（该对象被用作 Prototype，这是“四人帮”提出的另一种模式），这样该对象及其新兄弟将能够有效地共享 flyweight。

下面是该方法中的代码片段：

```
Entity * Entity::EntityFactory( int whichType )
{
    Entity * ptr;
    switch( whichType )
    {
        case SCAREDGHOST:
            ptr = new Ghost( SCAREDGHOST );
            break;

        case ACTIVEGHOST:
            ptr = new Ghost( ACTIVEGHOST );
            break;

        // ...
    }
    return ptr;
}
```

是不是很简单？调用该方法，并将一个导出类作为参数传递给它时，它将返回一个指针，该指针指向相应行为族的 **Entity** 子类。

```
Entity * newEntity = Entity::EntityFactory( ACTIVEGHOST );
```



附带光盘中的代码还实现了一个简单的窍门，用于装载文本文件中的关卡设置：实体的构造函数将关卡文件作为 **istream** 参数，并从该文件中读取实体所属类特有的参数。因此，工厂方法知道要创建的子类的内部情况。



### 1.8.9 在运行阶段选择策略

如果游戏中的行为类不多，或实体的行为很简单，无需使用脚本进行定义，则上面介绍的技术是管用的。然而，如果需要调整的内容很多，并需要使用脚本进行实验，该如何办呢？如果要求能够在运行阶段修改实体的策略，而不影响实体，又该如何办呢？这是策略（Strategy）模式的用武之地（我保证这是最后一个模式。我想是这样的。）

假设脚本编译器生成的是 C 语言代码。您需要将编译器生成的 C 函数与行为类（或各个实体）关联起来。这可以使用函数指针。

在 C++ 类中使用函数指针

要在运行阶段将一个方法插入到类中，最简单也是最好的办法是使用函数指针。在 C/C++ 中，函数指针是一个变量，就像其他指针一样，它存储的也是内存地址，只是被指向的对象是由无名特征标（即返回类型和参数列表）定义的函数。下面的代码声明了一个函数指针，它指向这样的函数：接受两个 Entity 对象作为参数，并返回一个 Boolean 值：

```
bool (*interactPtr) (Entity * source, Entity * target);
```

假设有一个特征标与此相同的函数，例如：

```
bool TypicalRabbitInteractions( Entity * source, Entity * target )
{
    ...
}
```

则可以将该函数赋给变量 `interactPtr`，然后通过解除该指针的引用来调用这个函数。因此下属两个代码片段是等效的：

```
Ok = TypicalRabbitInteractions( BasilTheBunny, BigBadWolf );
```

和

```
interactPtr = TypicalRabbitInteractions;
Ok = (*interactPtr) ( BasilTheBunny, BigBadWolf );
```

在类中使用函数指针要复杂些。其关键思想是，将脚本编译器生成的函数声明为类的友元，这样它便可以访问类的私有数据成员；同时将特殊的指针 `this` 作为第一个参数传递给该函数，`this` 表示当前对象。

```
class SomeEntity : public Entity
{
    // The function pointer
    void ( * friendptr ) ( Entity * me, Entity * target );

public:
    // Declare one or more strategy functions as friends
    friend void Strategy1( Entity * me, Entity * target );
```

```
...  
  
// The actual operation  
void HandleInteractions( Entity * target )  
{  
    (*friendptr) ( this, target );  
}  
};
```

这实际上相当于手工完成了调用函数时 C++ 编译器为您完成的工作：C++ 悄悄地将 `this` 作为参数传递给每一个方法。由于所有的现代编译器都将内联函数（inline function）调用，因此从性能上说，调用常规方法与调用编译后的脚本之间没有任何差别。

在运行阶段通过函数指针来选择策略也有助于减少层次结构中行为类的数目。极端情况下，可以使用单个 `Entity` 类替换整个类层次，该 `Entity` 类只包含策略元素的函数指针解引用。然而，这样做将使代码模糊不清，在某个地方完全不透明。因此，应谨慎而为。

最后，如果实体允许根据运行阶段的情况在多个不同的策略之间切换，则使用这种方案时，可以通过指针赋值来实现切换：快速、清晰、没有争议。

### 1.8.10 最后的注意事项

在简单情形下，本文介绍的技术甚至能够省却脚本编写。对于不需要功能全备的脚本语言，且/或负担不起与此相关的费用的小型项目，可以使用硬编码的策略代码、一个基于 GUI 的 SAMMy 编辑器和一个线性关卡描述文件（其中包含关于与每个实体相关的行为的键值元组）来实现：

```
<EntityName BasilTheBunny>  
<ExportedClass Rabbit>  
<StrategyVsEntity BigBadWolf Avoid>  
<HandleCollision BigBadWolf Die>  
...
```



附带光盘中包含多个组件类以及使用本文介绍的技术的范例。然而，要在项目中使用它们，需要做很多修改（例如，在 SAMMy 中添加 3D 模型）。

最后，用于装载 SAMMy 的文本文件格式以及代码中的其他对象被假设为脚本编译器、关卡编辑器或其他相关工具的输出。因此，它们的结构非常灵活，阅读起来可能很困难。如果您觉得阅读起来困难，请想想那些手工编写和编辑这些内容的可怜的作者吧。

### 1.8.11 参考文献

[Brown98] Brown, W.H., Malveau, R.C., McCormick III, H.W., Mowbray, T.J., *Anti Patterns: Refactoring Software, Architectures and Projects in Crisis*, Wiley Computer Publishing, 1998.

[GoF94] Gamma, E., Helm, R., Johnson, R., & Vlissides, J. (1994), *Design Patterns:*

*Elements of Reusable Object-Oriented Software*, Addison-Wesley, 1994.

[Rising98] Rising, L. ed., *The Patterns Handbook: Techniques, Strategies and Applications*, Cambridge University Press, 1998.



## 1.9 在 C++ 添加摒弃功能

Noel Llopis, Meyer/Glass Interactive  
nllopis@mgigames.com

在软件的整个生命周期中，函数接口可能需要修改、废弃或完全替换为新的接口。对于哪些被重用于多个项目或已用了多年的库和引擎，尤其是这样的。当函数与项目其他部分之间的接口被修改后，游戏（或和工具）可能无法通过编译。对于开发大型项目的小组而言，情况更糟，因为很多人可能会更频繁地修改接口。

### 1.9.1 可能的解决方案

对于这种情况，处理的方式有多种：

- **无为而治。**只要接口发生变化，继续工作之前所有的人都必须修改那些调用改动了的函数的代码。这对于一个人的团队而言是可行的，但对于大型团队而言，通常是不可接受的。
- **不修改接口函数。**这通常是不可能的，尤其是在瞬息万变的游戏行业。可能硬件变了，可能发行商需要新的东西，也可能最初的接口有问题。采用这种方式利大于弊，将导致函数或类的名称名不副实。
- **创建新的接口版本。**这种方式坚持不修改接口，而是创建新的接口来解决所有的问题。项目中将包含新的接口和原来的接口。DirectX 采用的就是这种方式。对于需要完全修改接口或不常更新接口的情况，这种方式管用；但对于经常需要更新或要做小的更新的情况，这种方法不太好。另外，采用这种方式时，需要维护当前接口的整个实现以及大量的旧接口，而这可能是一场恶梦。

显然，对于现代游戏开发，这些解决方案都不理想。我们必须采用其他方式来处理这种问题。

### 1.9.2 理想的解决方案

我们真正希望的是，能够编写新的接口函数，并将旧的接口函数保留一段时间。这样，小组的其他成员可以立即开始使用新函数。他们在有时间时对原来的代码进行修改，使之使用新的函数，经过一段时间后，当没有任何人使用旧的函数时，便可以将其删除。

这里的问题是，如何让每个人知道哪些函数已被修改以及他们被认为

将使用哪些函数。即使我们总是将这种信息告诉他们，但如果一切都能通过编译并正确运行，他们怎么会去记住这些呢？可以摒弃函数。编写新函数后，我们将函数标记为被摒弃的。这样，每当旧函数被使用时，编译器都将生成一条消息，指出调用的是一个被摒弃的函数，同时指出应使用哪个函数替换它。

### 1.9.3 使用和指定被摒弃的函数

Java 提供了一种特性，可以完成我们所需的功能。然而，当前的大多数商业游戏主要是使用 C++ 编写的，不幸的是，C++ 没有包含任何摒弃功能。本文余下的内容将介绍一个使用 C++ 实现的解决方案，用于将函数标记为被摒弃的。

首先来看一个如何使用该解决方案的例子。假设有一个名为 `FunctionA()` 的函数，每个人都需要使用它。不幸的是，几个月后，我们发现必须修改该函数的接口，因此我们编写了一个名为 `NewFunctionA()` 的函数。通过添加一行代码，我们可以将 `FunctionA()` 标记为被摒弃的。

```
int FunctionA ( void )
{
    DEPRECATE ( "FunctionA()", "NewFunctionA()" )
    // Implementation
}

int NewFunctionA ( void )
{
    // Implementation
}
```

这行代码 `DEPRECATE("FunctionA()", "NewFunctionA()")` 指出，`FunctionA()` 已被摒弃，并被 `NewFunctionA()` 所代替。

使用 `FunctionA()` 的用户无需做任何特殊的工作。每当用户使用 `FunctionA()`，并退出程序时，都将在调试窗口中发现下述消息：

```
*****
WARNING. You are using the following deprecated functions:
- Function FunctionA() called from 3 different places.
Instead use NewFunctionA().
*****
```

### 1.9.4 使用 C++ 实现摒弃功能



ON THE CD

一个 singleton 类 `[Gamma95]DeprecationMgr` 实现了所有这些功能。附带光盘中包含了这个类的完整源代码以及一个范例程序。在最简单的情况下，`DeprecationMgr` 只是维护了一个列表，其中包含所有已被摒弃的函数。每当 singleton 被释放时（在程序退出时自动进行），析构函数将在调试窗口中打印一个报表，指出当前会话使用了哪些已被摒弃的函数。

```
class DeprecationMgr
{
public:
    static DeprecationMgr * GetInstance ( void );
    ~DeprecationMgr ( void );

    bool AddDeprecatedFunction (const char * OldFunctionName,
                                const char * NewFunctionName,
                                unsigned int CalledFrom );

    // Rest of the declaration here
};
```

通常，我们无需直接处理这个类，因为宏 `DEPRECATE` 将为我们完成所有的工作。

```
#ifdef _DEBUG
#define DEPRECATE(a,b) { \
    void * fptr; \
    _asm { mov fptr, ebp } \
    DeprecationMgr::GetInstance()->AddDeprecatedFunction(a, b, fptr); \
}
#else
#define DEPRECATE(a,b)
#endif
```

除了前几行外，宏 `DEPRECATE` 所做的只是获取 `DeprecationMgr` 的一个实例，并将当前执行的函数加入到列表中。由于 `DeprecationMgr` 是一个 singleton，仅在 `GetInstance()` 函数被调用后才被实例化，因此，如果没有已被摒弃的函数，它将被创建，也不会再在程序执行完毕后打印任何报表。在内部，`DeprecationMgr` 为每个被摒弃的函数存储了一个小型结构，并通过一个 STL 映射集合按函数名对这些结构进行索引。仅当已摒弃的函数首次被调用时，才会该映射集合中插入一个新的条目。

`DeprecationMgr` 类还有另一项功能：记录在多少个不同的地方调用了每个被摒弃的函数。这很有用，因为这样我们一眼就可以看出，如果停止使用被摒弃的函数，需要在多少个地方进行修改。不幸的是，由于这直接使用了汇编，而它随平台而异，因此只能用于 x86 系列 CPU。`DEPRECATE` 宏的前两行获取 `EBP` 寄存器（通常可以从这里获得返回地址），并将其传递给 `AddDeprecatedFunction()`。然后，如果从同一个位置（例如在循环中）多次调用了函数，则只指出在一个地方调用了该函数。

这种获取返回地址的方法有一个潜在的问题。通常地址 `[EBP-4]` 包含的是当前函数的返回地址。然而，在有些情况下，编译器可能不会将寄存器 `EBP` 设置为期望的值。具体地说，在 VC++ 6.0 中，启用了编译器优化功能后，对于特别简单的函数可能会出现这种情况。在这种情况下，试图读取 `[EBP-4]` 要么返回的是一个不正确的值，要么将导致程序崩溃。在发布（release）模式下（此时优化功能通常被关闭），将不会有问题，因为该宏不做任何工作。然而，在调试模式下，有时候也会进行优化，因此在函数 `AddDeprecatedFunction()` 中，仅当 `[EBP-4]` 中包含的地址对当前进程而言是可读的，我们才试图去读取返回地址。这是通过使用异常处理或调用 Windows 特有的函数 `IsBadReadPtr()` 来实现的。这样，当优化功能启用时，报告的已摒弃的函数被调用的次数是错误的，但这样至少不会导致程序崩溃，而且摒弃管理器的



其他功能都将正确的发挥作用。

### 1.9.5 可改进的地方

---

还有一个重要的问题：摒弃警告是在运行阶段生成的，而不是在编译和链接时生成的。之所以必须这样做，是因为已摒弃的函数可能位于一个单独的库中，而不是被编译的代码中。只在运行阶段报告已摒弃的函数的惟一缺点是，程序可能仍在使用一个已摒弃的函数，但由于该函数调用的次数很少，因此没有被发现。直到已摒弃的函数被删除，导致编译器报错时，才会发现还在使用这个函数。

### 1.9.6 致谢

---

感谢 Dvid McKibbin 审校了此文，并指出编译器优化导致的问题及其解决方案。

### 1.9.7 参考文献

---

[Gamma95] Gamma, Eric, et al, *Design Patterns*, Addison-Wesley. 1995.

[Rose] Rose, John, "How and When to Deprecate APIs," available online at [java.sun.com/products/jdk/1.1/docs/guide/misc/deprecation/deprecation.html](http://java.sun.com/products/jdk/1.1/docs/guide/misc/deprecation/deprecation.html).



## 1.10 一个插入式调试内存管理器

Peter Dalton, Evans & Sutherland

pdalton@xmission.com

随着游戏编程变得越来越复杂，游戏的最低内存需求变得越来越高。为支持图形、音乐、视频、动画、模型、网络和人工智能，需要大量的资源，而当前的游戏必须高效地处理这些资源。随着项目不断增大，发生内存泄漏、超越内存边界以及分配过多内存的可能性也将增大。这正是内存管理器的用武之地。通过创建几个简单的内存管理例程，可以跟踪所有动态分配的内存，并指导程序优化地使用内存。

我们的目标是，通过报告内存泄漏、跟踪分配的内存被使用的比例，并提醒程序员不要边界违例，来确保内存被合理地使用。另外，我们还需确保到内存管理器的接口是无缝的，即不要求显式地调用函数或声明类。我们应该能够通过包含头文件，将这些代码毫不费力地插入到其他任何模块中。创建内存管理器的缺点是，管理器分配、释放以及查询内存以获得统计信息需要占用额外的时间。因此，最后创建游戏时，我们将不会启用该选项。为消除这些缺陷，我们只在调试时或符号 `ACTIVATE_MEMORY_MANAGER` 被定义时，才启用内存管理器。

### 1.10.1 内存管理器初步

内存管理器的核心思想是，对标准运算符 `new` 和 `delete` 进行重载，并使用 `#define` 来创建一些宏，这些宏可以插入我们自己的例程中。通过重载内存分配和释放例程，可以使用我们自己的内存跟踪模块替换标准例程。这些例程将记录请求内存分配的代码所在的文件和行号，并记录统计信息。

首先需要创建重载的 `new` 和 `delete` 运算符。正如前面指出的，我们将记录请求内存分配的代码所在的文件和行号。这些信息对于发现内存泄漏而言至关重要，因为通过这些信息，可以找到内存分配操作发生的位置。重载这些运算符的方式如下：

```
inline void*
    operator new(size_t size, const char *file, int line);
inline void*
    operator new[](size_t size, const char *file, int line);
inline void operator delete( void *address );
inline void operator delete[]( void *address );
```

需要注意的是,为确保正确运行,必须对运算符 `new` 和 `delete` 的标准版本和数组版本都进行重载。虽然这些声明不复杂,但我们现在需要解决的问题是:让所有将使用内存管理器的例程都无缝地将额外的参数传递给 `new` 运算符。为此,可以使用编译指令 `#define`。

```
#define new new(__FILE__,__LINE__)
#define delete setOwner(__FILE__,__LINE__),false ? setOwner("",0)
           : delete
#define malloc(sz) AllocateMemory(__FILE__,__LINE__,sz,MM_MALLOC)
#define calloc(num,sz)
           AllocateMemory(__FILE__,__LINE__,sz*num,MM_CALLOC)
#define realloc(ptr,sz) AllocateMemory(__FILE__,__LINE__,sz,
           MM_REALLOC, ptr )
#define free(sz) deAllocateMemory(__FILE__,__LINE__,sz,
           MM_FREE )
```

`#define new` 语句将把所有的 `new` 调用替换为新的 `new` 例程,后者接受的参数并非仅仅是需要分配的内存量,还包含文件和行号,以便跟踪内存的分配情况。Microsoft Visual C++ 编译器提供了一组预定义的宏,其中包括我们所需的 `__FILE__` 和 `__LINE__` [MSDN]。宏 `#define delete` 与 `#define new` 基本相同。在不引起语法问题的情况下,无法将额外的参数传递给重载的 `delete` 运算符,因此我们使用 `setOwner()` 方法来记录文件名和行号,供以后使用。另外,很重要的一点是,应有条件地创建这些宏,以避免多行宏中常见的问题 [Dalton01]。最后,出于完整性考虑,还必须使用我们自己的内存分配和释放例程来替换方法 `malloc()`、`calloc()`、`realloc()` 和 `free()` 等。



有关这些函数的实现,请参阅附带光盘。例程 `AllocateMemory()` 和 `deAllocateMemory()` 负责所有的内存分配和释放工作。它们还记录有关要求分配的信息,并根据要求对内存进行初始化或查询工作。然后,可以使用这些信息来生成统计信息,以分析程序的内存需求。

### 1.10.2 内存管理器的记录工作

提供了将标准内存分配例程替换为我们自己例程的框架后,便可以开始进行记录。正如本文开头指出的,我们将重点关注内存泄漏、边界违例的情况以及实际的内存需求。为记录所需的所有信息,必须首先选择一个数据结构来存储与内存分配相关的信息。为提高效率和速度,我们将使用一个链式哈希表。其中每个哈希表条目都将包含以下信息:

```
struct MemoryNode
{
    size_t      actualSize;
    size_t      reportedSize;
    void        *actualAddress;
    void        *reportedAddress;
    char        sourceFile[30];
    unsigned short sourceLine;
    unsigned short paddingSize;
```

```
char          options;  
long          predefinedBody  
ALLOC_TYPE    allocationType;  
MemoryNode    *next, *prev;  
};
```

该结构存储了分配给用户的内存数量，还存储了分配的内存块前后的补白（padding）内存量。我们还记录了分配类型，以防止分配与释放之间的不匹配。例如，如果分配内存时使用的是运算符 `new[]`，而释放内存时使用的是运算符 `delete`，而不是 `delete[]`，则可能由于没有调用对象的析构函数而导致内存泄漏。另外，在确保该结构具有最大灵活性的同时，要使其尽可能小。我们不想创建一个这样的内存管理器，即它使用的内存比被监视的应用程序使用的还多。

至此，我们拥有了确定程序中是否有内存泄漏所需的所有信息。通过在 `AllocateMemory()` 例程中创建一个 `MemoryNode`，并将其插入到哈希表中，可以记录分配的所有内存。然后，通过在 `deAllocateMemory()` 中删除 `MemoryNode`，可以确保哈希表只记录了当前分配的内存。如果在退出程序时，哈希表不为空，则说明发生了内存泄漏。在这种情况下，可以通过查询其中的 `MemoryNode`，将有关内存泄漏的细节告诉用户。正如前面指出的，在 `deAllocateMemory()` 例程中，我们还将检查用于分配内存的方法是否与用于释放内存的方法匹配。如果不匹配，我们将指出潜在的内存泄漏。

接下来收集有关边界违例的相关信息。当应用程序使用的内存超过分配给它的内存时，将发生边界违例。最容易出现这种情况的地方是访问数组的循环。例如，如果数组包含 10 个元素，而访问第 11 个元素，则将超越数组边界，重写或访问不属于该数组的信息。为防止这种问题，我们将在分配的内存前后提供补白。因此，如果一个例程请求分配 5 byte 的内存，`AllocateMemory()` 实际分配的内存将为  $5 + \text{sizeof}(\text{long}) * 2 * \text{paddingSize}$  (byte)。我们使用 `long` 来填充是因为它被定义为 32 bit 的整数。接下来，我们必须将补白初始化为一个预定义的值，如 `0xDEADC0DE`。这样，释放内存时，我们可以对补白进行检查，如果其值不是预定义的值，则说明发生了边界违例。在这种情况下，我们可以查询相应的 `MemoryNode`，并将边界违例情况告知用户。

需要收集的最后一项信息是程序的内存需求。我们希望知道分配了多少内存、其中被实际使用的内存有多少以及分配的最大内存量。为收集这些信息，我们需要另一个容器。下面的代码中只列出了类中与此相关的成员：

```
class MemoryManager  
{  
public:  
    unsigned int m_totalMemoryAllocations;  
    unsigned int m_totalMemoryAllocated;        // In bytes  
    unsigned int m_totalMemoryUsed;              // In bytes  
    unsigned int m_peakMemoryAllocation;  
};
```

在 `AllocateMemory()` 例程中，我们将能够更新 `MemoryManager` 中除变量 `m_totalMemoryUsed` 之外的所有信息。要确定分配的内存中实际被使用的有多少，需要采用

与用于确定边界违例的方法类似的技巧。通过在 `AllocateMemory()` 例程中将内存初始化为一个预定义的值，并在释放内存时查询其中的值，可以知道实际被使用的内存有多少。为获得好的结果，我们还是使用 `long` 值来初始化 32 bit 的内存边界，并使用一个预定义的值（如 `0xBAADC0DE`）来进行初始化。对于不处于 32 bit 边界内的其他字节，将被初始化为 `0xE` 或 `static_char<char>(0xBAADC0DE)`。虽然这种方法可能容易出错，因为没有预定义的值可用来将内存初始化为独特的，但初始化 32 bit 的边界将比初始化 1 byte 的边界好得多。

### 1.10.3 报告信息



有了所需的所有统计信息后，我们来解决如何将这些信息报告给用户的问题。附带光盘中的实现将所有的信息记录到一个日志文件中，一旦用户启用了该内存管理器，并运行程序，则在退出程序之前，将生成一个日志文件，其中包含所有的内存泄漏、边界违例和最终的统计报表。

最后一个问题是：我们如何知道程序何时将终止，以便输出日志信息？一种简单的解决方案是，要求程序员在程序终止之前显式地调用 `dumpLogReport()`。然而，这违背了创建无缝接口的原则。为在不使用显式函数调用的情况下，确定程序何时将终止，我们将使用一个静态类实例。其实现如下：

```
class Initialize
{ public: Initialize() { InitializeMemoryManager(); } };
static Initialize InitMemoryManager;

bool InitializeMemoryManager() {
    static bool hasBeenInitialized = false;
    if (s_manager) return true;
    else if (hasBeenInitialized) return false;
    else {
        s_manager = (MemoryManager*)malloc(sizeof(MemoryManager));
        s_manager->intialize();
        atexit( releaseMemoryManager );
        hasBeenIntialized = true;
        return true;
    }
}

void releaseMemoryManager() {
    NumAllocations = s_manager->m_numAllocations;
    s_manager->release(); // Releases the hash table and calls
    free( s_manager ); // the dumpLogReport() method
    s_manager = NULL;
}
```

我们要解决的问题是：确保内存管理器是第一个被创建的对象，同时也是最后一个被释放的对象。由于静态地定义的对象被处理的顺序，这很困难。例如，如果我们在创建内存管理器对象之前创建一个静态对象，而后者的构造函数动态地分配内存，则内存管理器将无法



跟踪这些内存。同样，如果我们使用`::atexit()`方法来调用一个负责释放内存的方法，则`::atexit()`方法被调用之前，内存管理器对象已被释放，因此指出的内存泄漏情况将是错误的。

为解决这些问题，需要作以下改进。首先，通过在内存管理器的头文件中创建 `InitMemoryManager` 对象，可确保它在任何静态对象声明之前被创建。在任何静态定义之前将内存管理器头文件包含进来时，情况也将如此。Microsoft 指出，静态对象被创建的顺序与出现顺序相同，而被释放的顺序则与此相反[MSDN]。其次，为确保内存管理器始终可用，我们将在 `AllocateMemory()`和 `DeallocateMemory()`中调用 `InitializeMemory()`，从而保证内存管理器处于活动状态。最后，为确保内存管理器是最后一个被释放的对象，我们将使用 `atexit()`方法。`::atexit()`调用传递给它的函数，但调用顺序同函数被传递的顺序相反[MSDN1]。因此，必须对内存管理器作出的惟一限制是，它必须是第一个调用`::exit`函数的方法。静态对象仍然可以使用`::atexit()`方法，只需确保内存管理器存在即可。例如，如果由于某种原因，`InitializeMemoryManager()`函数返回 `false`，则最后一个条件没有满足，因此日志文件将报告这种错误。

由于上述限制，使用 Microsoft Visual C++时，还有几点需要注意。VC++内部过程大量地使用了`::atexit()`方法，以便在关闭之前进行清理工作。例如，下面的代码将导致`::atexit()`被调用（要知道这一点，必须进行反汇编）：

```
void Foo() { static std::string s; }
```

虽然如果声明 `s` 之前内存管理器已处于活动状态，这不会导致任何问题，但还是有必要对此进行说明。虽然这个范例是针对 VC++的，但其他编译器可能与此不同或者包含其他在幕后调用`::atexit()`的方法。该解决方案的关键是：确保内存管理器被首先初始化。

#### 1.10.4 注意事项

要跟踪内存，需要占用内存和 CPU 时间，另外还有其他几个细节需要注意。首先，必须处理包含其他文件时可能导致的语法错误。在某些情况下，导致语法错误的原因可能是其他文件重新定义了运算符 `new` 和 `delete`。使用 STL 实现时，尤其容易出现这种情况。例如，如果我们包含 `MemoryManager.h`，然后再包含 `<map>`，则将导致各种类型的错误。要解决这种问题，我们将使用其他两个头文件：`new_on.h` 和 `new_off.h`。这些头文件将定义前面创建的 `new/delete` 宏或取消对它们的定义。采用这种方法的优点在于：不强迫用户遵守特定的 `#include` 顺序，因此更灵活；同时避免了处理预编译的头文件的复杂性。

```
#include "new_off.h"
#include <map>
#include <string>
#include <All other headers overloading the new/delete operators>
#include "new_on.h"

#include "MemoryManager.h" // Contains the Memory Manager Module
#include "Custom header files"
```

我们需要解决的另一个问题是：如何处理那些重新定义了运算符 `new` 和 `delete` 的库。例



如, MFC 有其处理运算符 `new` 和 `delete` 的系统[MSDN2]。因此, 我们让 MFC 类使用自己的内存管理器, 而让非 MFC 游戏代码使用我们的内存管理器。为此, 可以将 `#include "new_off.h"` 插入到 ClassWizard 创建的 `#ifdef` 后面。

```
#ifdef _DEBUG
#include "new_off.h"           // Turn off our memory manager
#define new DEBUG_NEW
#undef THIS_FILE
static char THIS_FILE[] = __FILE__;
#endif
```

这样, 我们能够保留 MFC 的内存管理器的优点, 如监视从 `CObject` 派生而来的类的内存泄漏情况, 同时继续为其他的代码提供一个内存管理器。

最后, 请记住必须正确实现 `delete` 运算符使用的 `setOwner()`。其实现不仅仅是记录文件和行号, 还必须创建一个堆栈实现。这是我们实现 `delete` 宏的方式决定的。例如, 对于下述代码:

```
File 1: line 1: class B { B() {a = new int;} ~B() {delete a;} };
File 2: line 1: B *objectB = new B;
File 2: line 2: delete objectB;
```

函数调用的次序如下:

```
1. new( objectB, File2, 1 );
2. new( a,      File1, 1 );
3. setOwner( File2, 2 );
4. setOwner( File1, 1 );
5. delete( a );
6. delete( objectB );
```

从上述代码可知, 当调用 `delete` 运算符来释放 `ObjectB` 时, 将无法获得有关文件和行号的信息, 除非我们使用堆栈来实现。虽然这个解决方案很简单, 但问题并非是显而易见的。

### 1.10.5 进一步的改进



在附带光盘提供的实现中, 对这里讨论的几点进行了改进。例如, 有一个选项, 让用户可以设置标记, 以执行更全面的内存检测。另外, 还有用于在内存分配和释放位置设置断点的选项, 以便能够查询程序的堆栈。这只是几项可能的改进, 还可以做其他改进, 如让程序检查给定的地址是否有效。在内存控制方面, 选项非常多。

### 1.10.6 参考文献

[Dalton01] Dalton, Peter, "Inline Functions versus Macros," *Game Programming Gems II*, Charles River Media. 2001.

[McConnell93] McConnell, Steve, *Code Complete*, Microsoft Press. 1993.

[MSDN1] Microsoft Developer Network Library, [http://msdn.microsoft.com/library/devprods/vs6/visualc/vclang/\\_pluslang\\_initializing\\_static\\_objects.htm](http://msdn.microsoft.com/library/devprods/vs6/visualc/vclang/_pluslang_initializing_static_objects.htm)

[MSDN2] Microsoft Developer Network Library, [http://msdn.microsoft.com/library/devprods/vs6/visualc/vccore/core\\_memory\\_management\\_with\\_mfc.3a\\_.overview.htm](http://msdn.microsoft.com/library/devprods/vs6/visualc/vccore/core_memory_management_with_mfc.3a_.overview.htm)

[Myers98] Myers, Scott, *Effective C++, Second Edition*, Addison-Wesley Longmont, Inc. 1998.



## 1.11 一个内置的游戏剖析模块

Jeff Evertt, Lithtech Inc.

jeff@evertt.com

本文介绍进行实时分析的低开销剖析（profiling）模块的架构和实现，它支持性能计数器组织，使得多用户能够协同工作。该模块是从游戏引擎的角度设计的，它的很多需求与游戏中的内容相关。撰写本文时，该模块已用于一个商业游戏引擎中。

几乎所有的人都认为，剖析游戏或引擎的性能是项重要的工作，只不过人们常常进行猜测或臆断，而不是创建一个能够收集真实数据的游戏系统。从长远的角度看，实现一个清晰的剖析系统是一项明智的投资。与其他任何事情一样，越早对剖析系统进行规划，实现它将越容易。

### 1.11.1 有关剖析的基本知识

剖析的基本机制很简单：分别在感兴趣的代码开头和末尾获取一个时间戳，然后将它们相减，得到的便是代码的运行时间。我们需要一个高精度计数器——Windows 多媒体定时器，其精度为 ms，因此很准确。如果平台为 PC 上的 Windows，则有两个高精度的 API 函数可用：QueryPerformanceCounter 和 QueryPerformanceFrequency。然而，由于这些函数的开销非常高，我们将创建自己的函数，这只需几行内嵌的汇编代码便可实现：

```
void CWin32PerfCounterMgr::GetPerfCounter(
    LARGE_INTEGER &iCounter) {
    DWORD dwLow, dwHigh;
    __asm {
        rdtsc
        mov dwLow, eax
        mov dwHigh, edx
    }
    iCounter.QuadPart = ((unsigned __int64)dwHigh << 32)
        | (unsigned __int64)dwLow; }
```

为将这一数字转换为秒，需要知道计数器的频率。这里，频率为每秒的 CPU 周期。我们可以在计数器启动时记录它的值——获取一个时间样本，然后暂停 500ms，并再次获取一个样本。注意，如果目标平台为游戏控制台，其中也有类似的计数器。

### 1.11.2 商用工具

对于性能调整而言，选择合适的工具至关重要。市面上有很多用于 PC 的时间检测工具，它们在应用程序运行时记录一个时间，然后脱机，以便能够逐个模块、逐个函数地查看特征数据。Intel VTune 和 Metrowerks Analysis Tool 都使用内置的 CPU 硬件计数器来生成处理后的游戏特征数据。VTune 的优点在于，能够通过指令排序或成对预测来调整汇编代码。

GPT (Intel Graphics Performance Toolkit) 提供了一些强大的实况分析工具。它在应用程序和 Direct3D/OpenGL 之间的层读取并查看数据。知道当前正被绘制的是什么很有用。有时候，改变游戏渲染的次序或方式将极大地影响性能。然而，GPT 是针对 DirectX 的特定版本编写的，因此其发行版本只能用于跟踪 DirectX。另外，获取重要的实况数据将降低应用程序的速度，因此使用 GPT 时，依靠获取的数据来推断性能特征将是很危险的。

图形卡的收集统计数据的驱动程序以及硬件计数器很有用。Nvidia 发布了特殊的驱动程序和一个实时数据查看应用程序，后者监视驱动程序的所有函数入口。如果图形驱动程序占用了大量的 CPU 时间，该应用程序让我们能够查看内部情况，并进一步分析。Intel 在其驱动程序和 i740 芯片硬件中提供了计数器，让我们能够在图形芯片级进行优化。有些游戏控制台也提供了这种功能。这种功能很有用，因为它是惟一种在如此低的级别上对性能进行分析的方式。然而，这要求对驱动程序和芯片的工作原理以及计数器的含义有深入的了解。

### 1.11.3 为何要自己开发模块

原因 1：基于帧的分析。在游戏中，帧与帧之间的相关性非常高，但在几秒钟之内又可能发生重大变化。请想象一下，一个 3D 射手（玩家）开始时面对一堵墙，然后沿长长的走廊奔跑，最后同 5 个 AI 驱动敌人展开血战。游戏引擎将经历很多可能不同的瓶颈，而只有逐帧进行分析，才能找出这些瓶颈。查看累积样值在整个时间段内的情况，将无法准确地了解发生的情况；而基于帧的分析则可以每次只针对一个问题。

原因 2：可以在任何时候、任何地方进行。在 PC 游戏开发周期的最后，有些人可能遇到性能问题，这种问题只在星期二在某种机器上出现。通常，有很多这种类型的问题。为调试这种问题，可能需要很多时间，导致游戏推迟发行。虽然这种问题是 PC 游戏中特有的，但控制台游戏却需要处理这样的问题，如在进入第三级后，发射导弹时，游戏速度变得非常慢。知道问题的原因后，解决起来通常很容易。如果可以在测试机器上显示一些计数器组，将可以迅速找出导致问题的罪魁祸首。

原因 3：可定制。现代的游戏引擎非常复杂。能够忽视引擎中的其他模块，而只关注正在运行的模块是一种强大的功能。另外，只有组织数据的工程师才能按其需求对数组进行组织。

#### 1.11.4 剖析模块 (Profiling module) 的需求

需求 1: 让用户能够快速、准确地对应用程序进行剖析。

需求 2: 开销非常低。如果帧时间中很大一部分是取样和显示结果引起的, 这实际上改变了系统中应用程序的行为。通常, 降低 CPU 的速度将隐藏图形卡引起的延迟。在有些情况下, 即使非常低的比例也将极大地改变游戏的性能, 因此, 当剖析模块被启用时, 其占用的 CPU 时间通常不应超过 5%; 被禁用时, 不应超过 1%。

需求 3: 允许多个用户在各自的系统上独立地工作, 而无需担心其他引擎模块。

需求 4: 在不需要时, 它不碍事。

#### 1.11.5 架构和实现

由一个性能计数器管理器 (IPerfCounterMan) 跟踪所有处于活动和非活动状态的计数器。计数器被划分为可被同时禁用和启用的组 (如模型渲染、背景渲染、AI、物理)。这样可以按照易于理解的分组方式, 将计数器表示为多个独立工作的群组。群组很有用的原因有两个: 可以快速确定是否需要某个计数器进行采样; 可以显示和隐藏整组计数器。我们将使用 4 个字符的编码 (FourCC) 来表示群组 ID; 对于计数器名称, 则使用全文字符串表示。

整个系统被组织成模块, 模块包含到其他模块的接口。基本组件是计数器, 它由群组 ID (FourCC) 和字符串名称标识。创建每个计数器时, 都赋予它一个整数 ID, 用于唯一地标识该计数器。通常, 游戏代码在初始化时创建计数器, 并在要剖析的代码前后分别启动和停止计数器。

该模块的基本功能单元接口如下:

```
class IPerfCounterMan {
public:
    // Add new counter (returns the ID, 0 is failure)
    int32      AddCounter(uint32 CounterGroup,
                        const char* szCounterName);

    // Forget your counter's ID? (Zero is failure)
    int32      GetCounterID(uint32 CounterGroup,
                        const char* szCounterName);

    // Delete the counter
    bool       DeleteCounter(uint32 CounterID);

    // Start and Stop a counter.
    void       StartCounter(uint32 CounterID);
    void       StopCounter(uint32 CounterID);

    // Draw the Counters onto the Screen (to be called once
    // per frame near the end of the scene)
    void       DrawCounters();
};
```

StopCounter 计算调用 StartCounter 和 StopCounter 之间的时间差, 获得代码运行的总时间。执行 DrawCounters 后, 所有运行的计数器都被清除。DrawCounters 还维护有一个设置到帧尾的最大值。我们假设游戏引擎有一个接受文本命令的调试控制台。启用和禁用计数器群组以及定制显示方式将很方便。

给计数器的显示方式提供尽可能多的配置很有帮助。我们不希望每帧都重新显示计数器(每 30 帧更新一次足够了), 但如果能根据调试的内容定制刷新时间将非常方便。另外, 同时显示当前的百分比和上一次刷新后的最大百分比也很有用。

用直方图显示结果是一种不错的方式。这样用户能够快速了解数据, 而且编写代码的工作也不难。另外如果能够在百分比和实际时间(单位为 ms)直接切换, 将时间和百分比显示为文本值, 自动缩放坐标轴, 将会很有帮助。不要频繁地改变坐标轴比例而不提醒用户, 这只会让用户感到迷惑。

### 1.11.6 实现的细节

到性能计数器管理器的接口应该灵活, 易于使用。剖析管理器的用户将发现, 与在初始化时一次性保存计数器相比, 随时调用 AddCounter(...)并提供完整的字符串, 以获得 ID 并启动计数器更容易。对于进行快速剖析而言, 提供这样的机制很有帮助。然而, 它的效率不高, 如果在一帧中调用这种机制多次, 占用的时间将急剧增加。另外, 提供一个可被放在函数开头的类是一种使用计数器的便利方式, 这个类的构造函数调用 StartCounter, 析构函数调用 StopCounter。

编写剖析管理器时, 最好提供某种#define 宏, 它能够将剖析模块完全删除。为使游戏的性能尽可能高, 首先应该删除剖析代码。我们需要为 AddCounter、StartCounter 和 StopCounter 提供宏, 使得当#ifdef 发现变化时, 便不将它们编译到游戏中。

另外, 最好通过颜色来实现可视化标记。当计数器被显示时, 如果每行的颜色不同, 则用户阅读起来将更容易。

### 1.11.7 分析数据

一定要剖析发行版本, 因为其瓶颈与调试版本可能不同。如果目标平台为 PC, 则最好选择两三种典型的系统配置(从低端到高端), 并对每一种配置进行剖析。在不同的系统配置之间, 瓶颈可能差别很大。

应对有性能问题的地方进行剖析, 同时应在典型的游戏进度中进行剖析。必须将问题分解, 每次重点关注一个方面, 并重点关注影响最大的地方。不能因为某个函数调用得最为频繁或占用了大部分的 CPU 时间, 就只将重点放在该函数上。我们常常只将周期时间同预期时间进行比较, 而实际的预期通常是根据经验得出的。

对剖析模块本身也应进行剖析。如果剖析模块的开销很高, 将改变游戏的行为。因此, 应在剖析模块的绘制例程前后也放置计数器。



### 1.11.8 有关实现的注意事项

---

这里描述的模块被实现成跨平台的。然而，其中的某些部分调用了依赖于平台的函数。实现时间戳查询和绘制函数的代码很可能依赖于平台，因此，最好在這些函数上添加一层抽象。这里描述的实现在绘制代码中使用了一组调试几何体和文本（其实现依赖于平台），这样它可以是独立于平台的。您可能需要编写一个创建 4 字符编码的宏，因为很多编译器不支持 4 字符编码。

这个系统也可用于对游戏服务器进行长时间的剖析，以发现其中的问题。所有的计数器的来源相同，因此很容易对数据进行过滤，并将其保存到磁盘中。



## 1.12 用于 Windows 游戏的线性编程模型

Javier F. Otaegui, Sabarasa Entertainment

Javier@sabarasa.com

以前，在 DOS 时代，我们基本上采用线性的方式来编写游戏。随后，创建工作从 DOS 转移到 DirectX 中，由于使用了 Windows 消息泵 (pump)，这是一个很大的跳跃。这种架构不能胜任游戏编程。本文介绍一种高效地封装消息泵的方式，它提供了一个线性编程模型，另外，作为一种不错的副作用，它允许通过“alt-tab”在应用程序间切换。此外，我们还将介绍如何正确恢复丢失的曲面。



如果您以前采用线性方式编写程序，将很容易理解本文介绍的方法的重要性。如果您一开始就是在 Windows 下编写游戏，会觉得，消息泵是游戏编程的固有环境，但您尝试过线性编程后，将再也不会愿意回过头来使用消息泵了。与大型状态机相比，它更为清晰，理解和调试起来也更容易。您可以节省大量设计、编程和调试时间，并会渴望能够以更为线性化的方式工作。

### 1.12.1 更新背景

现代游戏中经常包含某种 UpdateWorld 函数，它是消息泵中应用程序的核心，每当没有收到任何消息时都被调用。乍一看，UpdateWorld 函数的代码很容易编写：所有的应用程序变量、曲面和接口都已被初始化，只需更新并渲染它们即可。这可不是一项简单的任务，除非游戏中只有一个屏幕、没有场景切换、没有菜单，也没有选项。

问题在于：UpdateWorld 最终必须结束，并返回到消息泵，以便能够处理来自系统的消息。这样，我们将无法一直逗留在一个持续的 for 循环中。由于老式的 DOS 游戏无需频繁地返回到消息泵，以处理系统请求，因此可以以线性方式编写它们，而子程序可以连续地执行所需的循环、延迟或场景切换。我们只需在子程序中插入相应的代码即可。然而，使用消息泵后，我们必须时刻注意：每次循环后，都应返回。正如前面指出的，试图维护多个游戏屏幕时，每次循环后都需返回。

为解决这种问题，需要将应用程序中的每个子程序设计为一个有限状态机。每个子程序都将跟踪其内部状态，并根据该状态，调用其他几个子程序。被调用的每个子程序也是一个有限状态机，它执行完毕后（即没有

其他状态需要执行), 必须返回一个值, 让调用它的子程序知道, 现在可以继续处理自己的其他状态。当然, 每个子程序执行完毕后, 必须将自己的状态重置为 0, 让应用程序能够再次调用它。

假设有三四十个这样的子程序, 每个有几十种状态, 则我们面临的将是一个庞然大物。要调试或理解这些代码将非常困难。这种有限状态编程模型要比老式的线性 DOS 程序实现的模型复杂得多。

### 1.12.2 解决方案: 多线程 (Multithreading)

下面是一个简单的多线程模型, 它让游戏程序员能够避开消息泵及其有限状态编程模型。

Windows 支持多线程技术, 这意味着应用程序可以同时执行多个线程。这里的原理非常简单: 将消息泵放在一个线程中, 而将游戏放在另一个线程中。消息泵将留在最初的线程中, 因此我们能够将 UpdateWorld 函数从消息泵中取出, 将其恢复到最简单的形态 (线性编程模式)。现在, 我们只需在 doInit 函数中添加启动游戏线程所需的代码:

```
HANDLE hMainThread;           // Main Thread handle

static BOOL
doInit( ... )
{
    ... // Initialize DirectX and everything else

    DWORD tid;

    hMainThread=CreateThread( 0,
        0,
        &MainThread,
        0,
        0,
        &tid);

    return TRUE;
}
```

**MainThread** 被定义为:

```
DWORD WINAPI
MainThread( LPVOID arg1 )
{
    RunGame();
    PostMessage(hwnd, WM_CLOSE, 0, 0);
    return 0;
};
```

**MainThread** 将调用 RunGame 函数, 当它结束后, 我们只需发送一条 WM\_CLOSE 消息,

让消息泵线程结束执行。

### 1. 初始化代码

现在，我们必须决定将初始化代码（包括 DirectX 初始化代码）放在 `doInit` 函数中还是 `RunGame` 函数中。如果我们将所有响应 `WM_CLOSE` 消息的终止代码放在消息处理程序中的话，也许将其放在 `doInit` 函数中更合适。另一方面，我们也可以将所有的初始化代码放在 `RunGame` 函数中，这意味着我们将在新的线性游戏编程函数中直接处理代码中所有重要的部分。

### 2. “Alt-Tab” 问题

在游戏编程中，使游戏在 Windows 下是真正多任务的，是最难完成的工作之一。行为良好的应用程序必须能够正确地切换到其他应用程序。这意味着用户可以通过按下 `Alt-Tab` 来离开当前应用程序，有些游戏不允许这样，但我们将尝试实现这种功能。

可以使用标准的 `SuspendThread` 和 `ResumeThread` 函数，但这几乎不可能实现上述功能。我们将使用一个多线程的通信工具——事件。事件的工作原理类似于标记，可用于同步化不同的线程。我们的游戏线程将决定自己必须继续执行，还是必须等待到事件被设置。

程序启动时，我们必须创建一个手工重置的事件。该事件应在程序停用时被重置，并在程序再次激活时被设置。然后，在主循环中，我们只需等待事件被设置。

要创建事件，我们需要下述全局变量：

```
HANDLE task_wakeup_event;
```

要创建并设置事件，需要在初始化代码中包含以下代码：

```
task_wakeup_event =  
    CreateEvent(  
        NULL,           // No security attributes  
        TRUE,           // Manual Reset ON  
        FALSE,          // Initial state = Non signaled  
        NULL            // No name  
    );
```

在大多数游戏中都有一个这样的函数，即每当游戏需要渲染一个新的屏幕（通常，这时 DirectX 被请求刷新主缓存和备用缓存），都将在主循环中调用它。由于该函数被频繁调用，因此可以在其中使用下面的代码让线程等待事件进入空闲状态：

```
WaitForSingleObject( task_wakeup_event, INFINITE );
```

每当操作系统切换活动应用程序时，我们都必须将线程挂起。为此，我们必须有自己的 `WindowProc` 函数，并在收到 `APP_ACTIVATE` 消息时，检查应用程序是否处于活动状态。如果应用程序不处于活动状态，我们必须将游戏挂起。这需要执行下面的调用：

```
ResetEvent( task_wakeup_event );
```

然后将其恢复：

```
SetEvent( task_wakeup_event );
```

有了上述简单的实现后，当用户按下 Alt-Tab 后，游戏将暂停，腾出所有的处理器时间，让用户能够做其他工作。如果即使应用程序不处于活动状态，也必须继续更新背景，则我们可以将渲染管道（pipeline）挂起，并继续更新背景。该事件模型可用于任何数目的线程，只需为每个新线程插入新事件即可。

### 3. 处理失去的曲面

如果我们使用视频内存曲面，则当应用程序处于非活动状态时，将失去有关曲面的信息。我们现在面临的问题是，使用新的线性编程模型时，当程序被挂起时，可能正在执行一个子程序，这将导致曲面丢失。

对于这种问题，解决方案很多，其中之一是命令模式[GoF94]。不幸的是，这种方法会使代码晦涩难懂，而本文的主要目的是使代码更为清晰。我们可以使用一个由回调函数和 lpVoids 对组成的全局栈，需要重新装载曲面时，将调用回调函数。需要恢复曲面时，我们将调用 `callback_function(lpViod)`。参数 lpVoid 包含一些指针，这些指针指向我们需要的曲面，这样对于我们新的线性子程序而言，曲面将是本地的。

假设我们有一个名为 **Splash** 的子程序，它显示游戏中的飞溅屏幕——一个从文件中装载的曲面。如果当飞溅屏幕显示时，用户按下了 Alt-Tab，然后返回到游戏中，我们希望游戏再次显示飞溅屏幕（假设游戏不活动时，该曲面丢失了）。使用我们建议的方法时，必须这样做：

```
int LoadSplashGraphics( lpvoid Params )
{
    Surface *pMySurface;
    pMySurface = (Surface *) Params;

    // ...
    // (load the graphic from the file)

    return 1;
}

int Splash()
{
    Surface MySurface;

    // Push the function
    gReloadSurfacesStack.Push( &LoadSplashGraphics, &MySurface );

    // Do not forget to load graphics for the first time
    LoadSplashGraphics( &MySurface );

    // ... the subroutine functionality.

    // Pop the function
    gReloadSurfaceStack.Pop();
}
```

我们使用了一个栈，以便每个嵌套的子程序都能加入所有装载和生成曲面的代码。可以很容易地将该实现修改为另一个集合类，但由于其嵌套的功能，这是一个经典的面向栈的问题，因此栈更合适。

### 1.12.3 参考文献

---

[GoF94] Gamma, E., Helm, R., Johnson, R., & Vlissides, J. (1994), *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley. 1994.

Otaegui, Javier F., "Getting Rid of the Windows Message Pump," available online at [www.gamedev.net/reference/articles/article1249.asp](http://www.gamedev.net/reference/articles/article1249.asp).





## 1.13 栈缠绕

Bryon Hapgood, Kodiak Interactive  
bhapgood@kodiakgames.com

对于汇编程序员来说，栈缠绕（stack winding）是一种功能强大的技术，它让我们能够修改应用程序的栈，从而完成一些非常奇特的任务。只需做非常少的工作，我们就可以把这种技术扩展到 C/C++ 中。虽然使用机器语言编写每一行游戏代码的时代已一去不复返，但有时候，为使游戏更精致、速度更快，对这个神秘的领域做些探索还是值得的。

本文介绍一种特殊的栈缠绕，我称之为“暂时返回”。后续的范例将以这种最简单的形式为基础，最后获得 thunked 暂时返回。文中的代码已经使用 Microsoft MASM 和 Visual C++ 编译器测试过。我个人在很多 GameBoy Color、PC 和 Xbox 项目中使用过栈缠绕。

### 1.13.1 简单的 TempRet

顾名思义，栈缠绕是一种对栈进行修改，使之能够完成意想不到的工作的技术。术语栈缠绕来自这样一种思想：将值插入到已有的栈结构中，以修改其标准的和意料之中的行为。

程序清单 1.13.1 TempRet 例程

```
0  .586
1  .model flat
2  .data
3  buffer dd ?
4  file_handle dd ?
5  filesize dd ?
6  .code
7
8  _TempRetEg:
9
10     call    fn0
11     call    fn1
12     ;
13     ; before
14     ;
15     pop     edx
16     call    edx
17     ;
```

```
18      ; after
19      ;
20      call fn2
21      call fn3
22      ret
23
24 A:    call _TempRetEg
25      ret
26
27 end
```

程序清单 1.13.1 是栈缠绕的第一个基本组成部分：TempRet 例程。假设有一个函数（名为 MyFunc），它调用 \_TempRetEg，然后调用另外两个函数：fn0 和 fn1，最后是下述代码行：

```
pop edx
calledx
```

我们知道，CPU 处理第 24 行的汇编指令 CALL 时，将把下一行（第 25 行）的地址压入栈中，并执行 JUMP，跳转到第 8 行。第 15 行弹出栈中的地址，并将其存储到一个 CPU 寄存器中。现在 CALL 该地址。这将把第 20 行压入栈中，然后执行 JUMP，跳转到第 25 行。后者没有执行任何操作，而只是执行一个 CPU 返回，这将从栈中弹出一个地址，并跳转到该地址处。

继续执行 \_TempRetEg 中余下的代码，当该函数返回时，我们不返回到 MyFunc，而是返回到调用 MyFunc 的函数那里。这是一个有趣的小窍门，但它为何重要呢？当我们来看 fn0~fn3 时，其威力将显现出来。

假设 fn0 打开一个文件；fn1 分配一个缓冲区，并将该文件读取到内存中；fn2 释放分配的内存；而 fn3 关闭该文件。因此 MyFunc 无需关心释放的问题，它既不需要关闭文件，也无需释放内存。从功能上说，打开文件、将其读入到内存中、释放内存和关闭文件都在同一个代码块中。MyFunc 只需调用 \_TempRetEg，使用缓存，然后返回。

### 1.13.2 TempRet 链

将函数串联起来时，TempRet 范例不再管用。来看一个经典问题：DirectX 7 的初始化和释放。这通常需要很多步，但是以相反的次序释放 DX 组件是至关重要的，有时候，这可能非常复杂。

我们来扩展第一个范例，以演示这一点。

程序清单 1.13.2 将多个例程缠绕到栈上

```
0 .586
1 .model flat
2 .code
3
4 TempRet macro
5 pop edx
```

```
6 call    edx
7 TempRet endm
8
9 createWindow:
10     ; open the window
11     TempRet
12     ; close it
13     ret
14 setCooperativeLevel:
15     ; set to exclusive
16     TempRet
17     ; restore
18     ret
19 changeDisplayMode:
20     ; set 640×480 16 bpp
21     TempRet
22     ; restore
23     ret
24 createSurfaces:
25     ; create primary surface
26     ; get attached back
27     TempRet
28     ; release primary
29     ret
30 _SetupDX7:
31     call createWindow
32     call setCooperativeLevel
33     call changeDisplayMode
34     call createSurfaces
35     jmp _SomeUserRunFunc
36
37
38 end
```

通过连续执行多次 `TempRet`，我们实际上将 4 个例程缠绕到了栈上，这样当 `_SomeUserRunFunc` 返回时，将以相反的次序跳过 `createSurface`、`changeDisplayMode`、`setCooperativeLevel` 和 `createWindow`，到达 `TempRet` 后面的一行。

至此，我们使用的都是汇编语言，但要使用这种技术，并不一定非得编写汇编模块。在本文的最后一节，我们将介绍 Microsoft Visual C++ 中的两种机制：内嵌式汇编和裸体 (naked) 函数，它们能够帮助我们进行栈缠绕。

### 1.13.3 Thunking

现在，需要在 C/C++ 中实现前面讨论的思想。正如前面指出的，Visual C++ 一种便利的机制，可用于完成这样的工作，但其他编译器呢？如果不支持裸体函数，则必须使用汇编语言，因为栈结构的出现使问题变得复杂。这不是不可能，只是非常困难。

Thunking 是 Microsoft 推广的一种技术，用于忽略 (slipping) 两个代码段之间的一个代

码段。Thunk 之前，程序流程将沿代码执行。对于使用 C++ 实现栈缠绕模式而言，Thunk 是一种非常棒的方式。下面是一个范例，它完成了前面介绍的建立 DirectX 的任务。

程序清单 1.13.3: 使用 TempRet 的 Visual C++ 范例

```
#define TempRet\  
__asm{pop edx}\  
__asm{call edx}  
  
#define NAKED void __declspec(naked)  
#define JUMP __asm jmp  
#define RET __asm ret  
  
static NAKED createWindow(){  
    // open the window  
    TempRet  
    // close it  
    RET  
}  
  
static NAKED setCooperativeLevel(){  
    // set to exclusive  
    TempRet  
    // restore  
    RET  
}  
  
static NAKED changeDisplayMode(){  
    // set 640×480 16 bpp  
    TempRet  
    // restore  
    RET  
}  
  
static NAKED createSurfaces(){  
    // create primary surface  
    // get attached back  
    TempRet  
    // restore  
    RET  
}  
  
NAKED SetupDX7(){  
    createWindow();  
    setCooperativeLevel();  
    changeDisplayMode();  
    createSurfaces();  
    JUMP run  
}
```



### 1.13.4 递归

作为最后一个说明栈缠绕 (stack winding) 的强大功能的范例, 我们使用递归搜索 (bone animation key) 来解决一个经典的问题: 如何回滚 (roll back) 递归。在 C 语言中, 我们将不断地返回, 遍历整个栈, 直到到达栈顶。然而, 如果递归超过 100 层, 则这将占用一些时间。为解决这种问题, 有两个名为 SafeEnter 和 SafeExit 的实用函数。顺便说一句, 作为全局函数时, 这些代码在 C++ 对象中也管用。

程序清单 1.13.4 协助递归的 SafeEnter 和 SafeExit 函数

```
.586
.model flat
.code

public SafeEnter, SafeExit

; struct SAFE{
;   void*__reg[8];
;   void*__ret;
; }

; assembly for SafeEnter routine

_SafeEnter:

    pop edx ; return address
    mov eax,[esp] ; safe
    ;
    ;
    ;
    mov [eax].safe.__ret,edx
    mov [eax].safe.__ebx,ebx
    mov [eax].safe.__ebp,ebp
    mov [eax].safe.__esp,esp
    mov [eax].safe.__esi,esi
    mov [eax].safe.__edi,edi
    ;
    ;
    ;
    pop eax ; safe pointer
    pop edx ; call function
    push eax ; safe pointer
    mov ebp,eax
    call edx
    mov eax,ebp
    jmp sex
```

```
_SafeExit:

    pop edx ;    return
    pop eax ;    regs context
    ;
    ;
    ;
    mov edi,[eax].safe.__edi
    mov esi,[eax].safe.__esi
    mov esp,[eax].safe.__esp
    mov ebp,[eax].safe.__ebp
    mov ebx,[eax].safe.__ebx
    mov edx,[eax].safe.__ret
    mov eax,[eax].safe.__eax
    ;
    ;
    ;
    jmp edx

end
```

SafeEnter 函数将重要的 CPU 寄存器的内容复制到 SAFE 结构中，然后调用递归函数。就该函数而言，无需做其他的工作。找到所需的数据后，我们调用 SafeExit() 函数，并将前面创建的寄存器场境 (context) 传递给它。这样，将立刻返回到父函数。

现在，如果意想不到的情况发生，而搜索函数没有满足搜索条件，则函数只需以常规方式返回，沿调用链向上进行。

#### 程序清单 1.13.5 使用 SafeEnter 和 SafeExit 的递归范例

```
static void search(SAFE&safe,void*v){
    if(<meets_requirement>)
        SafeExit(safe);
    // do stuff
    search(safe,v);
    return;
}

int main(){
    SAFE safe;
    SafeEnter(
        safe,
        search,
        <some_pointer>)
    ;
}
```



## 1.14 自我修改的代码

Bryon Hapgood, Kodiak Interactive

bhapgood@kodiakgames.com

**自**我修改的代码也叫“RAM 代码”，它是一种引人入胜的技术，让程序能够在执行时修改自己的代码。这种技术已被用于从通用算法到神经网络（Neural network）等领域，并获得了令人惊奇的效果。在游戏中，这是一种功能强大的优化技术。最近，笔者在一个名为 Test Drive Cycles 的 GameBoy Color 游戏中使用这种技术来动态地对工艺图进行解压缩（速度为 60fps），对 14 个（而不是标准的 8 个）调色板进行解码，并实现了多级视差变换（parallax scrolling）。本文将介绍如何编写自我修改的应用程序。

### 1.14.1 RAM 代码的原理

RAM 代码蕴含的思想很简单，但要正确地使用需要耗用大量的时间。这种代码大部分是使用十六进制方式编写的，因此难以调试。我们来看一个非常简单的例子。我们要从 RAM 中装载一个 16 bit 变量中的指针。

```
get_hl:
    ld hl,ptr_var    ; Load HL register with the address ptr_var
    ld a,(hli)       ; Load A register with low byte
                    ; and increment HL
    ld h,(hl)        ; Load L register with high byte of ptr_var
    ld l,a           ; Save low byte into L
ret                ; Return
```

可以对上述代码进行改进，将其编写为：

```
get_hl:
    db $2a           ; ld hl,...
ptr_var
    dw $0C00         ; ...ptr_var
ret
```

从逻辑上说，这两个例程之间没有任何差别，但您能看出它们之间的其他差别吗？第二个范例将存储地址的变量作为中间值装载到 HL 中！换句话说，我们只是装载 HL，而不是出去装载一个地址。与存取主存储器相比，装载中间值的速度更快；另外，由于需要解码的字节数更少，因此代码的运行速度将快得多。保留寄存器时，我们可以进一步扩展这一思想。

我们预先将值写入到代码中，而不是压入和弹出所有的东西（这样做的成本很高）。例如，我们不这样编写代码：

```
get_hl:
    ld hl,ptr_var
    ld a,(hli)
    ld l,(hl)
    ld h,a
    ld a,(hl)
    push af ; Save A register
    ;
    ; do something with A
    ;
    pop af ; Restore A register
ret
```

而是将其优化为：

```
get_hl:
    db $2a      ; ld hl,...
ptr_var
    dw ptr_var ; ...ptr_var
    ld a,(hl)
    ld (var1),a
    ;
    ; do something with A
    ;
    db $2F      ; ld a,...
var1 db $00     ; ...saved register value
ret
```

就这个范例而言，这样做节省的成本并不多，但它说明了这种思想。

### 1.14.2 一个快速的 Bit Blitter

在很多游戏中，像素格式之间的转换至关重要，如从 16 bit(565)RGB 转换到 24 bit RGB。无论这种转换是在脱机工具或游戏中进行，这里介绍的例程都能满足需求。我们可以定义一个结构（名为 BITMAP），其中包含有关图像的信息。根据该结构，Blitter 能够使用 RAM 代码技术来创建一个执行缓冲区（execute-buffer）——一段使用 malloc 为其分配了内存的代码，其中包含的全部是汇编指令。

Blitter 接受两个例程，其中一个例程知道如何读取 16 bit(565)的 RGB 像素并将它们转换为 32 bit 的 RGB 值；另一个知道如何使用另一种格式书写它们。我们可以一次性地将这两个函数粘贴到一起（对宽度奇特的图像），或连续粘贴多次以展开（unroll）循环。下面的范例采用的是前一种方式。

我们来定义 BITMAP 结构和相关的枚举类型。

```
enum Format{
```

```

    RGB_3×8=0,
    RGB_565=4,
    RGB_555=8,
    RGB_4×8=12,
    RGB_1×8=16
};

struct BITMAP{
    void *pixels;
    u32 w,h,depth;
    TRIPLET *pal;
    Format pxf;
    u32 stride;
    u32 size;
~   BITMAP();
    BITMAP(int,int,int,Format,int n=1);
    void draw(int,int,BITMAP&,int,int,int,int);
    operator bool(){
        return v!=NULL;
    }
};

```

下面使用汇编语言定义该结构，这至关重要：

```

BITMAP struct
pixels    dd    ?
w         dd    ?
h         dd    ?
depth     dd    ?
pal        dd    ?
pxf        dd    ?
stride     dd    ?
_size     dd    ?
BITMAP ends

PF_BGR_3×8    =    00h
PF_BGR_565    =    04h
PF_BGR_555    =    08h
PF_BGR_4×8    =    0Ch
PF_BGR_1×8    =    10h

```

接下来，定义执行缓冲区：

```
execute_buffer db 128 dup(?)
```

为在 C++ 中使用这些代码，必须给成员函数 BITMAP::draw 取一个 mangled 名称。然后，需要编写初始化代码：

```

?draw@BITMAP@@@QAEXHHAU1@HHHH@Z:
    push ebp

```

```

lea ebp,[esp+8]      ; get arguments address
push ebx
push edi
push esi

mov edi,ecx          ; dst bitmap
mov esi,[ebp+8]      ; src bitmap
mov eax,[esi].bitmap.pxf

```

首先，我们必须确定是否需要转换。因此，我们检查 BITMAP 对象的两种像素格式是否相同：如果的格式相同，则检查它们的大小是否相同；如果大小也相同，则只需执行字符串复制操作即可；如果不是，但宽度相同，则仍可以进行字符串复制；如果宽度不同，则只能逐行地复制字符串。

```

mov edx,[edi].bitmap.pxf
cmp eax,edx
jne dislike
;
; like copy
;
mov ecx,[esi].bitmap._size
cmp ecx,[edi].bitmap._size
je k3
mov ecx,[edi].image.stride
mov edx,[esi].image.stride
cmp edx,ecx
jne @f
;
; same w different h
;
mov edx,[edi].image.h
mov eax,[esi].image.h
cmp eax,edx
jl k2
mov eax,edx
k2: mul ecx
mov ecx,eax
k3: mov esi,[esi].image.lfb
mov edi,[edi].image.lfb
shr ecx,2
rep movsd
jmp ou
@@:
;
; find smallest h --> ebx
;
mov eax,[edi].image.h
mov ebx,[esi].image.h
cmp ebx,eax
jl @f
mov ebx,eax

```

```

@@:      ;
        ;   calc strides
        ;
        add ebp,12
        mov eax,[ebp].rectangle.w
        mul [esi].image.depth; edx corrupts
        mov edx,[esi].image.stride
        sub ecx,eax
        sub edx,eax
        ;
        ;   calc offsets with intentional reg swap
        ;
        push eax
        push ecx
        push edx
        call calc_esdi
        pop  edx
        pop  eax
        pop  ecx
        ;
        ;   edx=dest pad
        ;
        shr ecx,2
        mov ebp,ecx
@@:      rep movsd
        lea edi,[edi+eax]
        lea esi,[esi+edx]
        mov ecx,ebp
        dec ebx
        jne @b
ou:      pop esi
        pop edi
        pop ebx
        pop ebp
        ret 1ch

```

如果两个位图的像素格式完全不同，则只能将每个像素从一种格式转换为另一种格式，代码如下。可以通过将循环展开来进一步改进这个例程——只需重复构建（build）步骤4次或更多次即可。

```

dislike: lea  eax,execute_buffer
        add  ebp,12
        push ou
        push eax
        push edi
        push esi
        push ebp
        mov  ebx,edi    ; destination image
        mov  edi,eax
        ;

```

```
; write "mov ebx,h"
;
mov al,0BDh
stosb
mov eax,[ebp].rectangle.h
stosd
;
; write "mov ecx,w"
;
mov al,0B9h
stosb
mov eax,[ebp].rectangle.w
stosd
;
; get read
;
mov edx,22
mov ebp,esi; source
mov eax,[ebp].image.pf
mov esi,rtbl_conv[eax]
lodsd
mov ecx,eax
add edx,ecx
rep movsb
;
; put write
;
mov eax,[ebp].image.pf
mov esi,wtbl_conv[eax]
lodsd
mov ecx,eax
add edx,ecx
rep movsb
;
; write tail
;
mov ecx,[esp]
push edx
sub dl,19
neg dl
shl edx,16
or edx,08D007549h
mov eax,edx
stosd ; start of exec_tail
mov eax,[ecx].rectangle.w; args
push eax
mov ecx,[ebp].image.stride; source
mul [ebp].image.depth
sub ecx,eax
jz @f
```



```

        mov     al,0B6h
        stosb
        mov     eax,ecx
        stosd
        jmp     pq
@@:      ;
        ;      modify outer branch
        ;
        dec     edi
        mov     eax,[esp+4]
        sub     eax,6
        mov     [esp+4],eax
pq:      pop     eax
        mul     [ebx].image.depth
        mov     ecx,[ebx].image.stride; dest
        sub     ecx,eax
        jz      @f
        mov     ax,0BF8Dh
        stosw
        mov     eax,ecx
        stosd
        pop     eax
        jmp     pr
@@:      pop     eax
        sub     eax,6
pr:      inc     al
        neg     al
        shl     eax,16
        or      eax,0C300754Dh
        stosd
        pop     ebp
        pop     esi
        pop     edi

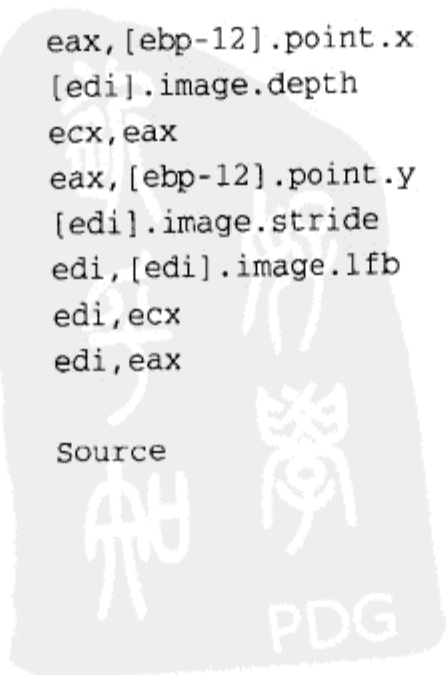
```

在该 **Blitter** 中,另一个重要的步骤是正确地计算源图像和目标图像的 x 轴和 y 轴偏移量,实现这项功能的例程如下:

```

calc_esdi:  ;
            ;      Destination
            ;
            mov     eax,[ebp-12].point.x      ; get dx
            mul     [edi].image.depth         ; multiply by d
            mov     ecx,eax                   ; store result
            mov     eax,[ebp-12].point.y      ; get dy
            mul     [edi].image.stride        ; multiple by stride
            mov     edi,[edi].image.lfb       ; get target pixels
            add     edi,ecx                    ; add x
            add     edi,eax                    ; add y
            ;
            ;      Source

```



```

;
mov     eax,[ebp].rectangle.x      ; get sx
mul     [esi].image.depth          ; multiply by d
mov     ecx,eax                   ; store result
mov     eax,[ebp].rectangle.y      ; get sy
mul     [esi].image.stride         ; multiple by stride
mov     edx,[esi].image.pal        ; palette info
mov     esi,[esi].image.lfb        ; get target pixels
add     esi,ecx                   ; add x
add     esi,eax                   ; add y
ret

```

为使该 RAM 代码管用，还需在 RAM 代码缓冲区的开头加入一些初始化代码，这些代码只是将要复制的扫描行数装载到 ECX 寄存器中。

```

exec_head    dd    5                ; (size)
              db    0B9h,000h,000h,000h,000h ; mov    ecx,0

```

下面是几个实际执行读写操作的例程（RC 和 WC）。第一个字节指出了每个例程有多少个字节的代码。

```

RC_BGR_1×8    dd    18                ; (size)
              db    033h,0C0h          ; xor    eax,eax
              db    0ACh                ; lodsb
              db    08Bh,0D8h          ; mov    ebx,eax
              db    003h,0C0h          ; add    eax,eax
              db    003h,0C3h          ; add    eax,ebx
              db    003h,0C2h          ; add    eax,edx
              db    08Bh,000h          ; mov    eax,[eax]
              db    025h,0FFh,0FFh,0FFh,000h ; and    eax,-1

RC_BGR_3×8    dd    7                ; (size)
              db    0ADh                ; lodsd
              db    025h,0FFh,0FFh,0FFh,000h ; and    eax,-1
              db    04Eh                ; dec    esi

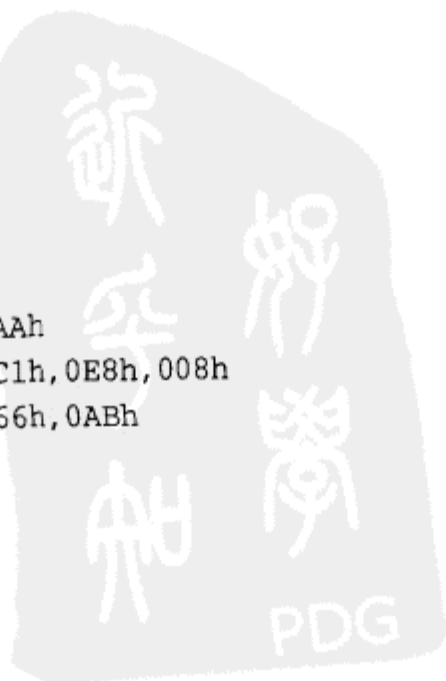
RC_BGR_4×8    dd    1                ; (size)
              db    0ADh                ; lodsd

RC_BGR_565:   dd    1                ; (size)
              lodsw

RC_BGR_555:   dd    1                ; (size)
              lodsw

WC_BGR_3×8    dd    6                ; (size)
              db    0AAh                ; stosb
              db    0C1h,0E8h,008h      ; shr    eax,8
              db    066h,0ABh          ; stosw

```



```

WC_BGR_555    dd    28                ;    (size)
               db    033h,0DBh         ;    xor ebx,ebx
               db    0C0h,0E8h,003h     ;    shr al,3
               db    0C0h,0ECh,003h     ;    shr ah,3
               db    08Ah,0DCh          ;    mov bl,ah
               db    066h,0C1h,0E3h,005h ;    shl bx,5
               db    00Ah,0D8h          ;    or bl,al
               db    0C1h,0E8h,013h     ;    shr eax,13h
               db    066h,0C1h,0E0h,00Ah ;    shl ax,0Ah
               db    066h,00Bh,0C3h     ;    or ax,bx
               db    066h,0ABh          ;    stosw

WC_BGR_565    dd    28                ;    (size)
               db    033h,0DBh         ;    xor ebx,ebx
               db    0C0h,0E8h,003h     ;    shr al,3
               db    0C0h,0ECh,002h     ;    shr ah,2
               db    08Ah,0DCh          ;    mov bl,ah
               db    066h,0C1h,0E3h,005h ;    shl bx,5
               db    00Ah,0D8h          ;    or bl,al
               db    0C1h,0E8h,013h     ;    shr eax,13h
               db    066h,0C1h,0ECh,00Bh ;    shl ax,0Bh
               db    066h,00Bh,0C3h     ;    or ax,bx
               db    066h,0ABh          ;    stosw

WC_BGR_4×8    dd    1                 ;    (size)
               db    0ABh              ;    stosd

```

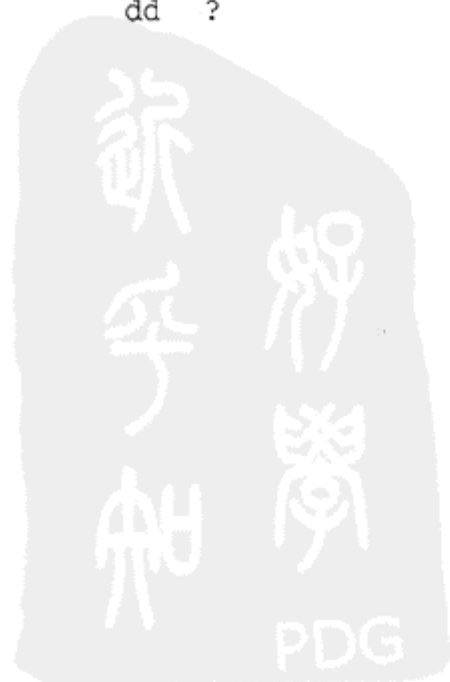
最后，需要一个表，它指出对于 `BITMAP::pf` 中的各种像素格式，应分别使用哪个例程。

```

rtbl_conv     dd    RC_BGR_3×8
               dd    RC_BGR_565
               dd    RC_BGR_555
               dd    RC_BGR_4×8
               dd    RC_BGR_1×8

wtbl_conv     dd    WC_BGR_3×8
               dd    WC_BGR_565
               dd    WC_BGR_555
               dd    WC_BGR_4×8
               dd    ?

```



## 1.15 使用资源文件来管理文件

Bruno Sousa, Fireworks Interactive

bsousa@fireworks-interactive.com

随着游戏越来越大（我想，最大的游戏应是 *Phantasmagoria*，有 7 张盘），必须对游戏数据进行组织。如果可执行文件所在的目录中有 10 个文件，是可以接受的；但如果有 1000 个，则是不可接受的。另外，如果目录结构超过 5 层以上，则处理起来将是极其痛苦的。由于游戏不可能类似于 Windows 资源管理器，因此必须找到一种简洁、快速的方式来存储和组织数据。这正是资源文件的用武之地。资源文件让我们能够采用有用的组织方式，将文件和目录封装到一个文件中；并能利用压缩、加密以及我们可能需要的其他任何特性。

### 1.15.1 何为资源文件

在日常工作中，我们使用过资源文件，如 Winzip、Windows 安装程序和备份程序。资源文件只不过是一种数据表示方式，这些数据通常来自多个文件，但被存储在一个文件中（参见程序清单 1.15.1）。通过使用目录，可使资源文件像硬盘驱动器的文件系统那样工作。

程序清单 1.15.1 资源文件的结构

```
Signature           = "SEALRFGNU" + '\0'
Version             = 1.0
Number of Files     = 58
Offset of First File = 19
```

```
[File 1]
[File 2]
[File 3]
[File .]
[File .]
[File .]
[File Number Of Files - 1]
[File Number Of Files]
```

资源文件中的每个块（以下将文件称作块（lump））都有自己的结构，结构的后面是数据（参见程序清单 1.15.2）。

### 程序清单 1.15.2 文件块的结构

```
File Size = 14,340
Filename = "/bmp/Bob.bmp" + '\0'
Flags = COMPRESSED
FlagsInfo = 0xF34A400B
[ Byte 1]
[ Byte 2]
[ Byte 3]
[ Byte .]
[ Byte .]
[ Byte .]
[ Byte File Size - 1]
[ Byte File Size]
```

### 1.15.2 设计

首先，需要给资源系统命名。这样，便可以使用该名称给每个部分指定一种特定的命名方案，将其与游戏的其他部分区分开来。我们将该系统称作“资源封装文件系统(Seal Resource File System)”，简称为 SRFS，并将“sl”作为类名前缀。

首先，我们需要一个资源文件头。从程序清单 1.15.1 可知，我们正在使系统简单明了。然而，这并不意味着其功能不强大，而是意味着它能提供所需的大部分特性，同时其语法和结构很容易理解。

我们的资源文件头提供了关于该系统的所有相关信息。游戏中使用了多种类型的文件，对于其中的每种文件，都有一个文件头，其中包含某些独特的信息，将其与其他类型的文件区分开来。SRFS 也不例外，因此其头中的第一项数据为文件特征标(file signature)。这通常是一个字符串，其中包含 5~10 个字符；它是必不可少的，这样我们才能确定该文件是一个有效的 Seal 资源文件。版本信息很简单——用于记录文件的版本，它也是必不可少的，原因很简单：如果我们决定对系统进行更新，以添加新特性或采用不同的方式对块进行排序，则需要核实使用的文件是否支持这些新特性。如果支持，则使用最新的代码；如果不支持，则应使用老代码——跨版本向后兼容是一个重要的设计问题，不应忘记这一点。头中的下一个字段用作特殊标记。在我们的文件系统的第一个版本中，没有使用这个字段，因此它必须始终为 NULL(0)。接下来是资源文件中包含的块数以及第一块的偏移量。该偏移量是必不可少的，这样可以返回到资源文件开头；另外它还可用来支持该系统的后续版本。可以在该头的后面加入有关后续版本的信息，而偏移量将指向第一个块。

现在来看块头，其中包含了检索数据所需的信息。首先是块的大小(单位为 byte)，然后是名称和目录，它们被存储为定长的、以 NULL 结尾的字符串。接下来是标记成员，它们指定了块的算法类型，如加密或压缩。然后，是关于算法的信息，这可以是加密的校验或压缩的字典信息(具体细节取决于算法)。最后，是以二进制方式存储的块信息。

我们的系统只有两个模块：资源文件模块和块(lump)模块。要使用块，需要从资源文件中将其装载到内存块中，并可能需要进行解密和解压缩，然后按常规方式存取内存块。游戏系统将所有功能都封装到资源文件模块中，甚至允许从该模块中直接存取块(lump)数据。

这种方式当然有其优点，但最大的缺点可能是：需要一次性将整个资源装载到内存中，除非只使用原始数据或使用复杂的算法来动态地将块数据解压缩或解密到内存中。这很困难，不在本文的讨论之列。

我们需要打开资源文件、读取头信息、打开各个块、读取块中信息以及获取块中数据的函数。这些函数将在“实现”一节介绍。

### 1.15.3 实现



附带光盘中的范例代码是使用 C++ 编写的，但在本文中，我们将使用伪代码，以便能够轻松地使用任何语言来实现。

#### 1. slClump 模块

slClump 模块类似于 C++ 中的文件流和其他语言的文件实现，因为我们可以将数据写入。不幸的是，由于 C++ 流的特性，更新资源文件中的块非常麻烦。我们不能将数据插入到流的中间，而只能替换它，因此无法修改资源文件。

```
DWORD      dwLumpSize;  
STRING      szLumpName;  
DWORD      dwLumpPosition;  
BYTE [dwLumpSize] abData;
```

变量 dwLumpSize 是一个双字 (32 bit)，它指定块的长度；szLumpName 是一个字符串，它描述了块的名称；dwLumpPosition 存储的是块的指针位置；而 abData 是一个字节数组，存储了块信息。

slClump 模块中的函数如下：

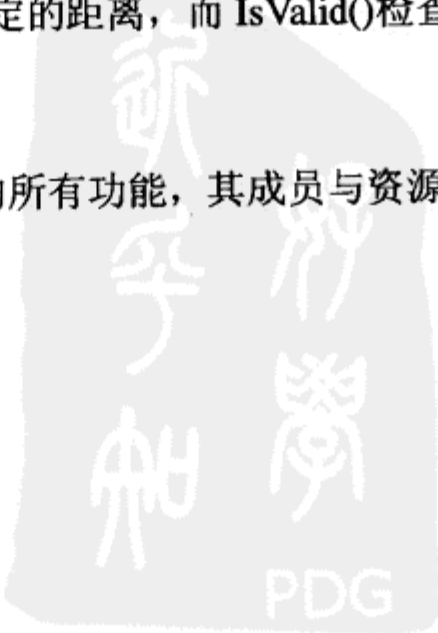
```
DWORD      GetLumpSize (void);  
STRING      GetLumpName (void);  
DWORD      Read (BYTE [dwReadSize] abBuffer, DWORD dwReadSize);  
DWORD      Write (BYTE [dwReadSize] abBuffer, DWORD dwWriteSize);  
DWORD      Seek (DWORD dwSeekPosition, DWORD dwSeekType);  
BOOLEAN     IsValid (void);
```

GetLumpSize() 获取块的长度，GetLumpName() 获取块的名称。Read() 将 dwReadSize 个字节读入到 sbBuffer 中，而 Write() 执行相反的操作，即将 sbBuffer 中的 dwReadSize 个字节写入。Seek() 移动块指针，使之远离某个位置指定的距离，而 IsValid() 检查某个块是否有效。

#### 2. slCResourceFile 模块

该模块具备装载资源文件中任何块所需的所有功能，其成员与资源文件头中的成员几乎完全相同。

```
DWORD      dwVersion;  
DWORD      dwFlags;  
DWORD      dwNumberOfLumps;
```



```

DWORD    dwOffset;
STRING    szCurrentDirectory;
FILE      fFile;

```

这些成员的用途已经在前面介绍过，这里简要地定义它们。`dwVersion` 是一个双字，它指定文件版本；`dwFlags` 是一个双字，它包含块的特殊标记；`dwNumberOfLumps` 是资源文件中的块数；`dwOffset` 指出了第一个块的位置（单位为 byte）；`szCurrentDirectory` 是我们所在的目录；而 `fFile` 是实际的 C++ 流。

下面介绍系统的核心部分——`slCResourceFile` 模块中的函数，我们将使用这些函数来存取各个块。

```

void      OpenLump (STRING szLumpName, slCLump inOutLump);
void      IsLumpValid (STRING szLumpName);
void      SetCurrentDirectory (STRING szDirectory);
STRING    GetCurrentDirectory (void);

```

其中每个函数都很简单。`IsLumpValid()` 检查名为 `szLumpName` 的文件是否位于资源文件中；`SetCurrentDirectory()` 将资源文件目录设置为 `szDirectory`。存取各个块时，将在块名前加上该目录名。`GetCurrentDirectory()` 返回当前目录。

下面来看 `Open` 函数。该函数打开资源文件中的一个块，下面是描述其算法逻辑的伪代码：

```

Check flags of Lump
  if Compressed
    OpenLumpCompressed (szLumpName, inOutLump)
  if Encrypted
    OpenLumpEncrypted (szLumpName, inOutLump)
  if Compressed and Encrypted
    OpenLumpCompressedEncrypted (szLumpName, inOutLump)
  else
    OpenLumpRaw (szLumpName, inOutLump)
end if

```



根据块的类型，将调用相应的函数来打开块。这样可确保设计优秀、代码简单。有关这些函数的源代码，请参阅附带光盘。

#### 1.15.4 有关实现的最后一些说明



对于有些用于打开文件或记录信息的支持函数，由于不能直接调用它们，因此前面没有列出。建议读者查看附带光盘中的源代码，其中做了良好的注释，因此很容易理解。压缩和加密算法是简单的 RLE 压缩和按位（bit-wise）加密，其实际的现实不在本文的探讨之列，必须单独研究。有关一些不受版权限制的（public domain）算法的信息，请参阅[Wostsit00]、[Wheeler00]和[Gillies98]。





### 1.15.5 结论

---

很容易将该系统进行改进或修改，以用于任意项目中。例如，可使其支持日期和时间有效性验证、添加防复制算法、校验和、数据池以及更好的压缩和加密算法等。

### 1.15.6 参考文献

---

[Hargrove98] Hargrove, Chris, "Code on the Cob 6," available online at [www.loonygames.com/content/1.11/cotc/](http://www.loonygames.com/content/1.11/cotc/), November 2~6, 1998.

[Towner00] Towner, Jesse, "Resource Files Explained," available online at [www.gamedev.net/reference/programming/features/resfiles/](http://www.gamedev.net/reference/programming/features/resfiles/), January 11, 2000.

[Wheeler00] Wheeler, David J, et al, "TEA, The Tiny Encryption Algorithm," available online at [www.cl.cam.ac.uk/ftp/users/djw3/tea.ps](http://www.cl.cam.ac.uk/ftp/users/djw3/tea.ps).

[Wotsit00] Wotsit.org, "The Programmer's File Format Collection: Archive Files," available online at [www.wotsit.org](http://www.wotsit.org), 1996~2000.

[Gillies98] Gillies, David A. G., "The Tiny Encryption Algorithm," available online at <http://vader.brad.ac.uk/tea/tea.shtml>, 1995~1998.



## 1.16 游戏输入的记录和重放

Bruce Dawson, Humongous Entertainment  
bruced@humongous.com

**18**世纪的数学和物理学家 Marquis Laplace 指出,如果能够知道宇宙中每个粒子的位置、方向和速度,将可以通过一个公式来详细预测未来和过去[Reese80]。这就是决定论。

混沌理论、海森堡的测不准原理和量子物理中的随机性表明,决定论是错误的。然而,在简化的游戏领域中,Laplace 的决定论却是管用的。

如果您仔细记录影响游戏的每一项内容,便可以重放记录,让昔日重现。

### 1.16.1 记录输入有何用途

记录游戏输入的用途比很多人想象的要多:重现难得一见的 bug、重放有趣的游戏、度量优化效果或创建游戏电影。

#### 1. 重现 bug

计算机程序是确定的,并完全可以预测,然而我们经常听到有人说,遇到了难以重现的 bug,因此难以排除。如果说程序是确定性的,为何 bug 难以重现呢?

有时候,罪魁祸首是硬件或操作系统。线程切换时间的选择和硬件驱动器并非完全一致,因此,代码中的竞态条件可能导致罕见的崩溃。然而,这种罕见的崩溃常常是由于不常见的用户输入引起的。在这种情况下,至少从理论上说,bug 是可以重现的——如果我们能够重现用户的输入序列的话。

对测试过程进行录像将有助于找出这样的 bug,然而,如果时机至关重要,这种方式将毫无帮助。为何不利用计算机的可预测性,让计算机程序记录所有的输入,并在必要时重放呢?

这里的关键是,要通过记录输入来找出 bug,必须确保即使游戏崩溃时,也能将输入记录下来。在 Win32 上,这非常容易。通过创建一个结构化异常处理(structured exception handling)程序[Dawson99],可以在游戏崩溃时,保存输入缓冲区。

如果在游戏引擎中加入一个选项,用于“快进”游戏输入过程(只渲染一部分帧),则可以快速重现崩溃过程。如果我们添加另一个选项,用于渲染游戏更新循环(update loop)处的帧,则可以轻松地重现崩溃时的情景。

要重现 bug，需要记录并重放所有的用户输入，包括与菜单屏幕的交互。菜单代码也可能存在 bug。

## 2. 重放有趣的的游戏

记录游戏输入最常见的用途是，让玩家能够记录有趣的的游戏。这些记录被用来演示如何玩、制作教程、测试新计算机的硬件或共享游戏。

为记录游戏，让用户以后能够重放，最重要的一点是，必须始终进行记录。期望用户在游戏开头决定是否要记录游戏是不现实的；应在游戏结束后，询问用户是否要永久性存储记录下来的游戏。

## 3. 评估优化效果

对优化而言，最重要的是测量优化前后的性能。如果不这样做，很可能出现这样的情况，即优化后的代码实际上运行速度更慢。

度量游戏性能是项复杂的工作，因为它变化多端。多边形数目、纹理设置、路径搜索的复杂性以及场景中的对象数目等都会影响帧速。测量单个帧的时间是毫无意义的，而快速浏览也是不科学的。

重放游戏输入是一种不错的解决方案。如果重放多次，并记录有关游戏性能的详细信息，便可以将每次修改后的情况绘制成图表，从而了解修改后的效果以及哪些地方还需改进。记录普通情况和最糟情况下的帧速以及帧速的连贯性将更容易，也更有意义。

并非总是可以通过重放游戏输入来检测优化效果，因为您的修改可能影响游戏的行为——修改后，帧速非常高可能是由于玩家进入了一个密室。因此，通过重放游戏输入来检测优化效果时，必须记录重要的游戏状态，并查看重放时有何变化。

## 4. 创建游戏电影

要创建演示片，可以将 VCR 与具备视频功能的图形卡相连，然后再玩游戏。然而，结果可能不太完美。VCR、视频编码器和游戏不断变化的帧速将导致影像模糊、抖动。

通过记录游戏输入，录制有趣的的游戏，然后播放它将是小菜一碟。只需对引擎做些微的修改，便能够告诉引擎：已进入有趣的部分，并从实时播放模式切换到电影录制模式。在这种模式下，引擎每秒能够精确地渲染 60 帧，并将每一帧记录到磁盘中。帧速可能降低到每秒两帧，但没有关系，因为录制的输入将完美地重放。

## 5. 实现多玩家

很多游戏（*X-Wing vs. TIE Fighter* 和《帝国时代》）都将输入记录和重放技术用于其网络模型中[Lincroft99]。它们只是传输玩家的输入，而不传输玩家的状态信息。这种技术尤其适用于有成千上万个个体的策略游戏。

### 1.16.2 原理

从理论上说，记录游戏输入很简单，实践上同样也可能很简单。然而，还是有些微妙之

处，如果您不小心，它们将引发问题。

### 1. 使游戏是可预测的

为使游戏输入记录和重放管用，游戏必须是可预测的。换句话说，游戏必须不受任何不可预测或未知因素的影响。例如，如果游戏可能受任务切换时机的影响，则它将是不可预测的。

很多游戏使用可变的更新和渲染循环交替，并以设置的频率记录输入。渲染一帧后，根据需要运行更新循环若干次，以处理累积下来的输入。

这种模型表明，每次渲染一帧后，游戏更新循环的运行次数是不确定的；然而，这并不会导致游戏本身是不可预测的。如果您要查找渲染函数（`render`）中的 bug，则可能需要知道渲染循环和更新循环之间的交替细节。但在其他情况下，这是毫无关系的。有必要记录每帧发生了多少次更新，但重放时可以忽略这些信息，除非您要查找渲染函数中的 bug。

然而，如果渲染函数修改了游戏的状态，且更新循环和渲染函数之间的交替是不固定的，则游戏将是不可预测的，因此记录输入将不管用。一个这样的例子是，渲染函数和更新循环使用的是同一个随机数生成器。*Total Annihilation* 也有一个这样的例子，在该游戏中，仅在渲染场景后，才更新“战争烟雾”（Fog-of-War(FOW)）。这种优化是有道理的，因为它降低了这种成本高昂的操作的频率。虽然它确保了玩家只能看到准确的烟雾，但却导致游戏是不可预测的。AI 使用的战争烟雾与渲染函数相同，因此渲染函数的调用时机将影响游戏的进程。

可能导致游戏不可预测的其他因素有：未被初始化的局部变量和不总是返回结果的函数。对于这两种情况，游戏的行为将取决于栈中的内容。这属于代码中的 bug，必须找出它们。

一个可能导致不可预测性的、难以处理的问题是声音重放，这是因为声音硬件以异步的方式处理它们。声音硬件的微小变化都将导致音效延迟结束。即使变化非常小，但如果正好位于两帧之间，而游戏将等待声音结束的话，这将影响游戏的行为。

对于很多游戏而言，这不是什么问题，因为无需将游戏与声音同步。如果要求同步，有一种很有效的解决方案——近似。启动音效时，计算样本将播放多长时间（样本数除以频率）。然后等待指定的时间过去，而不是等待声音结束。结果几乎相同，但将保持高度一致。

### 2. 初始状态

您还需确保游戏的起始状态是已知的——无论是重新开始新游戏还是装载保存的游戏。这通常是自动发生的。然而，每当您重新编译或修改数据时，都将稍微改变初始状态。幸运的是，对代码和数据所做的很多修改并不会影响游戏展开的方式。例如，如果您修改纹理的大小，帧速可能发生变化，但行为不会——如果游戏是可预测的话。如果对纹理大小的修改导致其他所有的内存块被分配到不同的地方，这也不会有任何影响——只要代码中不存在任何内存重写 bug。

一种修改代码和数据可能影响游戏行为的情况是，修改了生物或墙壁的位置或稍微调整了某种事件发生的概率。细微的修改可能不会带来变化，但会破坏可预测性。

在浮点计算中，结果可能发生意料之外的变化。编译优化后的版本时，编译器生成的代

码计算得到的结果可能与优化前的版本不同,有时候,这种差别可能会有影响。您可以使用 Visual C++ 中的优化器设置 Improve Float Consistency 来将这种问题降低到最小,但浮点误差是一个无法避免的问题,您只能提防。

### 3. 随机数

随机数可用于确定性游戏中,但有几点需要注意。

可以使用随机数的原因是,rand()并非真正随机的。rand()是使用一种简单算法实现的——通常是一种线性叠合方法,这种算法通过了很多随机数测试,但完全是可再现的。这被称为伪随机数生成器。只要您使用一个一致的种子初始化 rand(),便可得到一致的结果。如果希望游戏每次的随机情况不同,则可以根据时间来为 rand()选择种子,但将该种子记录下来。这样,如果要重现游戏,则可以重用该种子。

rand()的问题之一是,它生成单个随机数流(stream)。如果渲染代码和游戏更新代码都使用的是 rand(),且每秒渲染的帧数是不固定的,则该随机数生成器将很快变得不确定。因此,游戏更新循环和渲染函数应使用不同的生成器来获得随机数,这至关重要。

rand()的另一个问题是,其行为是不可移植的。也就是说,其行为可能随平台而异。

如果您保存游戏,然后继续玩,随后想重新装载保存的游戏并重放以后的输入,rand()的第三个问题便将显现出来。为使这是可预测的,必须将该随机数生成器重置为保存游戏时的状态。困难在于,无法实现这样的重置。C 和 C++ 标准没有对 rand()生成的随机数与为重置其状态而需要传递给 srand()的数字之间的关系做任何说明。例如,Visual C++ 在内部维护一个 32 bit 的随机数,但 rand()只返回其中的 15 bit,使得无法重新设置种子。

基于这三个问题,可以得出这样的结论:不要使用 rand();而应创建可移植、可重新启动的随机数对象。您可以让渲染循环使用一个随机数对象,而让游戏更新循环使用另一个。



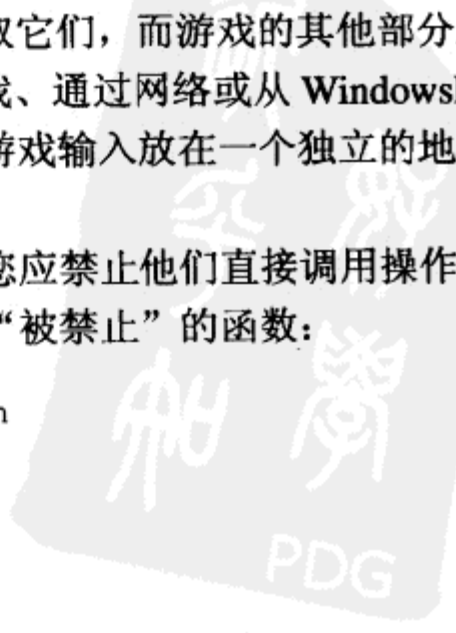
实现自己的随机数对象时,请不要自己发明一种随机数算法。随机数生成器极其微妙,您不太可能自己发明一种优秀的算法。请参阅附带光盘中的范例代码、Web 资源[Coddington]或 Knuth 编写的图书[Knuth81]。

### 4. 输入

恢复游戏的初始状态后,应确保能够记录并重放将影响游戏的所有输入。如果游戏更新循环直接调用操作系统函数来获取用户输入——如调用 Win32 函数 GetKeyState(VK\_SHIFT)来判断 Shift 键是否被按下,则这项工作很难完成。相反,所有的输入都需通过一个输入系统。该系统可以记录每帧开始时,所有输入设备的状态,并在游戏更新循环请求时,提供这些信息。输入系统可以将这些信息存储到磁盘中或从磁盘中读取它们,而游戏的其他部分对此一无所知。输入系统能够读取直接输入的数据、从保存的游戏、通过网络或从 WindowsProc 读取数据,而更新循环对这种差别一无所知。另外,最好将游戏输入放在一个独立的地方,这样游戏代码将更清晰,可移植性也更高。

程序员习惯于违反那些没有明确实施的规则,因此,您应禁止他们直接调用操作系统的输入函数。可以使用下面的技术来防止程序员无意间使用“被禁止”的函数:

```
#define GetKeyState Please do not use this function
```





```
#define GetAsyncKeyState Please do not use this function either  
...
```

对于多玩家游戏而言，另一项重要输入来自网络。如果希望能够重放游戏，则必须记录来自网络的数据和玩家的输入流。这样，即使没有网络连接，您也能够重放游戏。网络数据流可能非常大——通过 56KB 调制解调器运行的游戏，每小时可能收到几兆的数据。虽然这种大型数据流可能使记录更加难以进行，但还没有达到引发问题的程度。记录该数据流的好处数不胜数，但成本却非常低。

游戏可能使用的最后一项输入是“时间”。您可能希望在某个时间发生特定的事件，这些时间必须是游戏时间，而不是实际时间，这至关重要。每当游戏需要知道时间时，都应通过游戏引擎获取当前的游戏时间——除非用于剖析。与其他输入函数一样，最好使用预处理器来防止有人无意间编写了这样的代码，即调用了 `timeGetTime()` 或其他操作系统的时间函数。

记录这个游戏过程的输入是个不错的主意，这样，您将能够通过重放输入来查找游戏中任何地方的 bug，即使是 bug 位于游戏开始前的菜单中。然而，出于很多目的，您可能想在游戏进程中单独存储输入记录，以便能够单独重放它。

### 1.16.3 测试输入记录功能

---

游戏输入记录功能应该可用于任何编写优秀的游戏。即使游戏是多玩家的，只要记录机器收到所有的输入，也能够重现游戏。

然而，如果重放游戏的结果同以前不一致，则难以找出其中的原因。为找出问题，一种不错的方法是：只记录部分游戏状态——这可能是所有游戏实体的生命值和位置。这样，重放时，便可以检查不同的地方，并在差别显现之前发现它们。

### 1.16.4 结论

---

游戏输入的记录和重放是游戏引擎的重要组成部分，它能带来很多好处。如果一开始就进行规划，则很容易加入这种功能，并将使游戏引擎的设计更完美，更灵活。下面是一些应遵循的规则：

- 让所有的游戏输入（包括键盘、鼠标、游戏杆、网络和时间）都经过一个输入系统，以确保一致性，并能够记录和保存所有的输入。应自始至终地记录输入，并在游戏崩溃时或玩家请求时，永久性地存储这些输入。
- 提防代码中的浮点数优化或 bug，它们可能导致优化后版本的行为不同或不可预测。
- 不要使用 `rand()`，而应使用随机数对象。
- 决不要在渲染函数中修改游戏的状态。
- 存储游戏进程中的一些游戏状态，以便能够自动检测到不一致的情况。这有助于发现竞态条件、无意间修改的代码或 bug。



附带光盘中的范例代码中包含一个输入系统和一个随机数类。

### 1.16.5 参考文献

---

[Reese80] Reese, W.L., *Dictionary of Philosophy and Religion*. Humanities Press, Inc. 1980. p. 127.

[Knuth81] Knuth, Donald, *The Art of Computer Programming, Second Edition, Volume 2, Seminumerical Algorithms*.

[Coddington] Coddington, Paul, "Random Number Generators," available online at [www.npac.syr.edu/users/paulc/lectures/montecarlo/node98.html](http://www.npac.syr.edu/users/paulc/lectures/montecarlo/node98.html).

[Dawson99] Dawson, Bruce, "Structured Exception Handling," *Game Developer* magazine (Jan 1999): pp. 52~54.

[Lincroft99] Lincroft, Peter, "The Internet Sucks: What I Learned Coding X-Wing vs. TIE Fighter," 1999 Game Developers Conference Proceedings, Miller Freeman 621~630.





## 1.17 一个灵活的文本分析系统

James Boer, Lithtech Inc.

jimb@lithtech.com



几乎每个现代游戏都需要某种文本分析程序。本文（和光盘中的范例代码）展示了一个功能强大、易于使用的文本分析系统，该系统可用于处理任何类型的文件格式。

在表示数据方面，文本文件有很多优点：

- 它们易于使用任何标准文本编辑器进行读取和编辑。二进制数据通常要求创建、调试和维护定制的工具。
- 它们灵活——可使用同一个分析程序进行简单的变量赋值或编写更复杂的脚本。
- 它们可以在代码和数据之间共享常量（后面将更详细地介绍）。

不幸的是，文本文件也有些缺点：

- 不同于大多数二进制格式，文本必须首先进行标记（tokenize）和解释，因此装载速度更慢。
- 在空间方面，存储文本的效率不高，它浪费磁盘空间并降低文件装载速度。

由于很多游戏参数只需在开发阶段进行调整，因此可以在开发阶段使用基于文本的格式，而在产品中使用更佳的二进制格式。这样便融合了两方面的优点：文本数据易于使用；二进制数据的装载速度快。本文的后面将介绍如何将文本文件编译为二进制格式。

### 1.17.1 分析系统

该分析程序支持以下功能：

- 对基本数据类型（关键字、运算符、变量、字符串、整数、浮点数、布尔值和 GUID）的本机（native）支持；
- 能够识别用户定义的关键字和运算符；
- 支持 C 语言和 C++ 风格的注释（多行注释和单行注释）；
- 能够读写编译后的二进制文件；
- 调试功能——发生错误时，指出源代码文件和行号；
- 可使用 #include 包含了的文件；
- 支持 #define 宏替换。

上述大部分内容的含义都是不言自明的，但就文本分析程序而言，对

`#include` 和 `#define` 的支持好像令人费解。我们将讨论这些特性将如何极大地简化脚本，并提供另外一种机制来防止脚本与代码不同步。

### 1.17.2 宏、头文件和预处理技术

按处理 C 或 C++ 代码的方式对数据文件进行预处理有一些好处。解释这种概念的最佳方式是来看一个简单的例子。假设我们要使用一个脚本文件来创建大量独特的对象，脚本文件将提供正确的初始化每个对象所需的数据，并创建供代码使用的独特句柄。这种脚本的内容可能像下面这样：

```
CreateFoo(1)  { Data = 10 }
CreateFoo(2)  { Data = 20 }
CreateFoo(3)  { Data = 30 }
CreateBar(4)  { Foo = 1 }
```

假设关键字 `CreateFoo()` 将触发代码创建一个 `Foo` 对象，则内存中将有 3 个 `Foo` 对象，其中每个对象都有其独特的成员数据（这是由脚本创建的）。另外，假设我们将通过句柄来引用这些对象，则可以在代码中将 1、2 和 3 作为独特的句柄来存取这些对象。注意，在我们的例子中，脚本也可以使用这些数字式句柄。`Bar` 类需要一个有效的 `Foo` 对象作为其成员数据，因此当我们创建第一个 `Bar` 对象时，使用指向第一个 `Foo` 对象的引用。

创建数百个句柄后，将难以记录每个句柄的值。每当在脚本中添加一个对象时，程序员都必须在代码中修改同样的值。这无法避免程序员无意间引用了错误的脚本对象的情况发生。在 C 和 C++ 中，通过使用头文件解决了这种问题。在头文件中，可以声明变量或其他常用的元素，供多个源文件共享。如果将文本脚本视为一个源文件，则类似于 C 语言中的预处理器的优点将立刻显现出来。我们对前面的范例进行修改，以使用头文件而不是数字。

- Header File -

```
// ObjHandles.h
// Define all our object handles
#define SmallFoo 1
#define MediumFoo 2
#define LargeFoo 3
#define SmallBar 4
#define FooTypeX 10
#define FooTypeY 20
#define FooTypeZ 30
```

- Script File -

```
//
// Directs the parser to scan the header file
#include "ObjHandles.h"

CreateFoo(SmallFoo)      { Data = FooTypeX }
CreateFoo(MediumFoo)    { Data = FooTypeY }
```



```
CreateFoo(LargeFoo)      { Data = FooTypeZ }  
CreateBar(SmallBar)      { Foo = SmallFoo }
```

修改后，除了不再包含数字从而更容易阅读和理解外，文本脚本和源代码还可以共享头文件，因此它们不可能会不同步。

由于我们使用`#define`实现了一种简单的预处理替换，因此分析和使用更复杂的宏时，只比这多一个步骤。现在，通过识别通用的基于参数的宏，可通过替换参数来简化复杂的脚本操作。宏便于使用的原因还有一个：由于仅当宏被实际使用时，才会被编译到代码中（类似于脚本的原始形态），因此我们可以创建定制的基于脚本的宏，而不会破坏头文件的 C++ 兼容性。

注意，虽然我们处理了宏和`#define`，但分析程序无法识别其他命令，如`#ifdef`、`#ifndef`和`#endif`等。

### 1.17.3 该分析系统的结构

我们的分析系统有 5 个类：`Parser`、`Token`、`TokenList`、`TokenFile` 和 `Macro`。`Macro` 是一个用于 `Parser` 类内部的助手（helper）类，因此我们只需了解在 `Parser` 内是如何使用它的。`TokenFile` 是一个可选类，用于读写标准标记列表中的二进制标记。余下的是分析系统的核心（`Parser`、`Token` 和 `TokenList`）。由于 `Token` 是分析程序生成的基本构件，我们首先来介绍它。

#### 1. Token 类

该分析系统的基本数据类型是 `Token` 类。这个类可存储 8 种不同的数据类型：关键字、运算符、变量、字符串、整数、实数、布尔值和 `GUID`。其中，关键字、运算符、变量、字符串都是用 C 字符串表示的，因此它们之间的惟一差别是语义。整数、实数和布尔值是用 `signed int`、`double` 和 `bool` 表示的。对于大多数用途而言，这应该足够了。`GUID`（全局惟一的标识符）也是用本机数据类型表示的，因为使用一种可保证惟一性的数据类型将比较方便。

`Token` 类包含类型字段和共用体，后者可用于存储多种不同类型的数据。这样用一个类便可以表示所有的基本数据类型。存取数据时，首先检查 `token` 处理的数据类型，然后调用相应的 `Get()` 函数。断言（`assert`）确保了不会进行不合适的数据存取。

每种数据类型都在分析程序中发挥了一定的作用。理解它们的工作原理至关重要，这样可避免脚本错误。总体而言，类型定义类似于 C++ 中的定义。所有的关键字和标记都是区分大小写的。

#### 2. 关键字

关键字是专门定义的单词，它们被存储在分析程序中。有两个预定义的关键字：`include` 和 `define`。用户定义的关键字主要用于协助标记的词法分析（扫描阶段后）。

#### 3. 运算符

运算符通常是由一或两个字符组成的符号，如赋值运算符和逗号。运算符比较独特，它们类似于空白，因为它们可将其他数据类型分开。因此，在扫描例程中，运算符的优先级总

是最高的，这意味着运算符中的符号不能用于关键字或变量名中。因此，不应将数字或字母用作运算符的组成部分。对于本分析系统中的运算符还有另一种限制：由于使用了搜索方法，大于单个字符的运算符都必须是由更小的运算符组成的。当没有被空白或其他标记隔开时，大型符号总是优先于小型符号。

#### 4. 变量

变量是任何没有出现在关键字列表中的、基于字符的标记。

#### 5. 字符串

字符串必须用双引号括起。本分析程序支持最长可达 1024 个字符的字符串（在分析程序中，缓冲区常量是可以调节的），但不支持多行字符串。

#### 6. 整数

本分析程序能够识别正数和负数，并将它们存储为有符号的整型值。它还能够识别十六进制数（前缀为 0x）。不执行范围检查。

#### 7. 浮点数

浮点数使用 `double` 值表示。本分析程序能够识别任何带小数点的数字，但不懂科学计数法。对于浮点数，也不进行范围检查。

#### 8. 布尔值

布尔值是使用 C++ 中的 `bool` 类型表示的，而 `true` 和 `false` 是内置的关键字。与 C++ 一样，这些值也是区分大小写的。

#### 9. GUID

通过使用宏扩展代码，可以支持 GUID，而不需要额外做太多的工作。注意，除非使用 `ProcessMacros()` 来扩展宏，否则 GUID 将仍然是一系列的基本类型。该函数将在后面介绍。

#### 10. TokenList 类

`TokenList` 类是从标准的 STL `list of Tokens` 派生而来的，因此具备后者的所有特性，同时它还新增了两个特性。`TokenList` 类允许查看标记所在的文件和所处位置的行号，这有助于调试。使用编译开关可以删除这种功能。

#### 11. Parser 类

这是分析功能的核心。我们首先创建一个 `Parser` 对象，并调用 `Create()` 函数。注意，所有的函数都返回一个 `bool` 值——如果成功，则为 `true`；如果失败，则为 `false`。接下来，我们必须保留除默认外的其他运算符和关键字，供文本分析使用。

然后，便进入实际的分析阶段。分析阶段有三步，由 3 个函数进行处理。将该功能细分让用户能够更好地控制分析过程。通常，对于简单的分析任务，不需要处理包含文件和宏替

换。第 1 步是读取文件，并将文本转换为 `TokenList`，这是使用函数 `ProcessSource()` 完成的；接下来，`ProcessHeader()` 查找其中的头文件，然后对其分析，并将其替换为原来的源代码；第 3 个函数 `ProcessMacros()` 完成简单和复杂的 C 风格宏替换。这是一种非常强大的功能，它很有用，尤其对脚本语言而言。

来看看整个处理过程。出于简化和清晰的目的，没有进行任何错误检查。

```
// We need a Parser and TokenList object to start
TokenList toklist;
Parser parser;

// Create the parser and reserve some more keywords and tokens
parser.Create();
parser.ReserveKeyword("special_keyword");
parser.ReserveOperator("[");
parser.ReserveOperator("]");

// Now parse the file, any includes, and process macros
parser.ProcessSource("data/scripts/somescrpt.txt", &toklist);
parser.ProcessHeaders(&toklist);
parser.ProcessMacros(&toklist);
```

## 12. TokenFile 类

由于处理和分析人类能阅读的文本文件很慢，因此在发行代码中必须使用一种效率更高的文件格式。`TokenFile` 类可以将处理后的标记列表转换为二进制形式，这可以避免分析文本文件多次、进行 `#include` 搜索和宏替换等。基于字符的值，如关键字、运算符和变量被存储在查找表中。所有的数字值都是以二进制格式存储的，这样降低了所需的空间，也提高了效率。通常，二进制格式的装载速度为文本格式的 5~10 倍。

`TokenFile` 类的用法也很简单。`Write()` 函数接受一个 `TokenList` 对象作为参数，并使用指定的文件名或输出流创建二进制格式。另外，这个类在存储文件时，既可以区分大小写，也可以不区分。如果脚本中同时包含变量 `Foo` 和 `foo`，则关闭大小写区分功能后，在二进制格式中它们将被合并在一起，从而可以进一步节省空间。默认情况下，这种功能被关闭。

文件的读取工作是由 `Read()` 函数完成的，该函数的代码与下面类似：

```
TokenFile tf;

// Write a file to disk
tf.Write("somefile.pcs", &toklist);

// Or read it
tf.Read("somefile.pcs", &toklist);
```

### 1.17.4 小结

在最简单的情况下，文本文件的处理非常简单，只需几行代码。然而，对于其他更复杂的情况，最好编写一个功能齐全的文本分析系统，其健壮性和灵活性取决于这种工作的需求。

## 1.18 一个通用的调节器

Lasse Staff Jensen, Funcom

lasse@funcom.com

在游戏开发过程中，最常见的任务之一是：对变量进行调整，直到游戏取得我们希望的平衡。本文将介绍一个易于使用的“调节器(tweaker)”接口及其实现所涉及的设计问题。

### 1.18.1 需求分析

通用调节器接口的主要目标之一是：尽可能地透明和易于使用。在这里，用户是要调节其变量的程序员。需要强调的其他需求包括：使用的内存量、能够调节变量而不会增加太多的开销以及调节变量的速度（因为在有些情况下，该调节器也将被用于发行版本中）。

对这些需求做进一步分解后，可以知道实现需要完成的任务：

- 它应该对编码人员透明，这意味着要调整的变量不应包含其他的数据和/或功能，而且使用这些变量时，无需了解调节器。
- 它应该易于使用，这意味着定义将被调整的变量时，用户需要编写的代码不应超过 10 行，同时只需编写两行代码，便可以调整并取得变量。

### 1.18.2 实现

#### 1. 设计方案

图 1.18.1 是一个 UML 类图，本文余下的内容将从下到上依次介绍这些类。该设计方案中的基本要素是类型信息和调节层次结构。

#### 2. 类型信息

我们将使用模板具体化(specialization)来提供类型信息，这些信息能够以统一的方式存储。首先是基类 TypeID\_c，它使用一个虚函数定义了类型信息类的接口，该函数返回一个包含类型名称的字符串。

```
class TypeID_c{
public:
    virtual const char* GetTypeName() const { return "Unknown"; }
};
```







接下来，我们创建一个模板类，它可用于获得指定变量的类型。在这个类中，我们添加了一个成员函数，该函数返回一个指向 `TypeID_c` 实例的指针，可以直接检查该指针存储的地址。

```
template <class T>
class Identifier_c {
public:
    static const TypeID_c* const GetType();
};
```

声明上述类后，我们将使用模板具体化来定义每种类型。`TypeID_c` 的每个子类都将是一个 `singleton`，而指向该实例的指针将作为类型标识符。出于简化的目的，它们都将被声明为静态成员，从而使作用域为全局。我们可以通过接受 `GetIdentification` 方法返回的指针，来确保从其他函数调用时，实际的实例是存在的。对于 `float` 值，完整的实现如下：

```
class floatID_c : public TypeID_c {
public:
    virtual const char* GetTypeName() const { return "float"; }
    static TypeID_c* const GetIdentification();
};

template <>
class Identifier_c<float> {
public:
    static const TypeID_c* const GetType() {
        return floatID_c::GetIdentification();
    }
};

TypeID_c* const floatID_c::GetIdentification() {
    static floatID_c cInstance;
    return &cInstance;
}
```

要使用这些类来获得类型信息，我们只需存储基类指针：

```
float vMyFloat;
...
const TypeID_c* const pcType = TweakableBase_c::GetTypeID( vMyFloat );
```

其中，`TweakableBase_c`（后面将更详细地介绍）有一个模板成员，它调用正确的 `Identifier_c` 具体化。然后，我们便可以检查指针的地址：

```
if( Identifier_c<float>::GetType() == pcType ) {
    // We have a float!
}
```



在附带光盘的代码中，有两个用于定义用户类型的宏。因此要支持新的数据类型，只需在头文件中调用 `DECLARE_DATA_TYPE`，并在实现文件中调用 `DEFINE_DATA_TYPE`，然后重新编译即可（另外，如果不想支持范围检查，可能还需要调用宏 `DUMMY_OPERATORS()`）。

### 3. TweakableBase\_c

我们已有了一种简单、明了的存储类型信息的方式，接下来创建一个基类，用于包含指向可调节变量的指针。这个类还包含一个模板成员，用于获取前面提到的类型信息。由于需求之一是确保内存开销尽可能小，因此将使用 `RTTI` 来检查被存储到内存中的可调节变量的类型。然后，我们添加一个虚函数来确保这个类的多态性，该函数返回被存储的类型信息（如果没有，则返回 `NULL`）。这个类的实现如下：

```
class TweakableBase_c {
public:

    TweakableBase_c( void* i_pData ) : m_pData( i_pData ) {}
    ~TweakableBase_c() { /*NOP*/; }

    virtual constTypeID_c* const GetStoredType() const { return NULL; }

    template <class T>
    static constTypeID_c* const GetTypeID( const T& i_cValue ) {
        return Identifier_c<T>::GetType();
    }

protected:
    void* m_pData;

}; // TweakableBase_c
```

有了基类后，便可以创建包含其他数据（如类型信息、对范围检查的限制、用于回调函数的指针以及我们要关联到各种可调节变量的其他数据）的子类，同时确保使用的内存量最少。下面是一个可调节类的代码：

```
template <class T>
class TweakableType_c : public TweakableBase_c {
public:

    TweakableType_c( T* i_pxData, constTypeID_c* i_pcType ) :
        TweakableBase_c( reinterpret_cast<void*>( i_pxData ) ),
        m_pcType( i_pcType ) { /*NOP*/; }

    constTypeID_c* const GetDataType() const { return m_pcType; }
    virtual constTypeID_c* const GetStoredType() const { return m_pcType; }

private:
    constTypeID_c* const m_pcType;
```

```
}; // TweakableType_c
```

上述代码的优点在于，子类被实现为模板，虽然基类不是使用模板定义的。这样，我们便可以传递指向实际数据类型的指针，对接口隐藏到 void 的强制类型转换。

#### 4. Tweaker\_c

至此，我们具备了创建调节器类所需的所有部件。Tweaker\_c 类将存储所有的可调节变量，并提供给用户对这些被存储的值进行调节的功能。我们将使用一个 STL 映射表来存储所有指向可调节变量的指针，并将可调节变量的名称作为键值。所有的功能都是由简单模板成员提供的，下面是 TweakValue 成员：

```
template<class Value_x>
TweakError_e TweakValue( const std::string& i_cID, const Value_x& i_xValue ) {
    TweakableBase_c* pcTweakable;
    iTweakableMap_t iSearchResult = m_cTweakable_map.find( i_cID );
    if( iSearchResult == m_cTweakable_map.end() ) { return e_UNKNOWN_KEY; }
    pcTweakable = (*iSearchResult).second;

#ifdef _DEBUG
    TweakableType_c<Value_x>* pcType;
    if( pcType = dynamic_cast< TweakableType_c<Value_x>* >( pcTweakable ) ) {
        assert( pcTweakable->GetTypeID( i_xValue ) == pcType->GetDataType() );
    }
#endif

    TweakableTypeRange_c<Value_x>* pcTypeRange;
    if( pcTypeRange = dynamic_cast< TweakableTypeRange_c<Value_x>* >( pcTweakable ) )
    {
        assert( pcTweakable->GetTypeID( i_xValue ) == pcTypeRange->GetDataType() );

        if( i_xValue < pcTypeRange->GetMin() ) { return e_MIN_EXCEEDED; }
        if( i_xValue > pcTypeRange->GetMax() ) { return e_MAX_EXCEEDED; }
    }

    *(reinterpret_cast<Value_x*>( pcTweakable->m_pData ) ) = i_xValue;
    return e_OK;
} // TweakValue
```

由于成员是一个模板，因此我们可以直接将它转换为给定的值，从而完全隐藏了 void 转换。注意，如果用户决定不存储类型信息，则很容易强制我们做出糟糕的事情，因为我们无法检查 reinterpret\_cast 的来源！

#### 5. TweakerInstanceDB\_c

为支持对可调节变量进行分组以及存储一个变量的多个实例的功能，我们使用一个实例数据库来存储不同的调节器。其实现很简单——一个 STL multimap，其中存储了不同调节器的所有实例；以及一个由这些 multimap 组成的 STL 映射表，其中类别被用作键值。

### 1.18.3 使用

下面根据需求对该实现进行测试，看我们是否实现了预期的目标。要将变量定义为可调节的，需要创建一个调节器，并将其加入到可调节变量实例数据库中。

```
Tweaker_c* pcTweaker = TweakerInstanceDB_c::AddTweaker( "Landscape",
    TWEAKER_CREATE_ID( this ), "Graphics" );
```

上述代码为 **Landscape** 类创建了一个调节器，并将其加入到 **Graphic** 类别中。**TWEAKER\_CREATE\_ID** 宏将 **this** 指针作为参数，并确保 **Landscape** 类的每个实例都有一个独特的 ID。然后，我们将每个变量加入到这个调节器（和其他调节器）中：

```
pcTweaker->AddTweakable( &m_vShadowmapScaleTop, "Shadowmap scale", 0.0F, 68.0F );
```

上述代码加入了一个变量，将其取值范围设置为[0, 68]，并将其命名为“**Shadowmap scale**”。注意，由于 **AddTweakable** 方法的模板特性，我们必须传递正确类型的参数（例如，使用 **0.0F** 而不是 **0**）。要将变量定义为可调节的，需要两行代码；同时，对使用该变量的用户完全隐藏了这一点。

要调节该变量，我们需要的只是名称、数据类型和实例。通常，我们拥有指向调节器实例本身的指针，但在 GUI 代码中，常常可以这样做：

```
TweakerInstanceDB_c::iConstCategoryMap_t iCategory =
    TweakerInstanceDB_c::GetCategory( "Graphics" );
Tweaker_c* pcTweaker =
    GetTweaker( iCategory->second, "Landscape", TWEAKER_CREATE_ID( pcLandscape ) );
```

在上述代码中，我们首先取得所有属于 **Graphic** 类别的实例，然后搜索 **Landscape** 类的实例（我们假设指针 **pcLandscape** 指向所需的实例）。要修改特定变量的值很简单：

```
Tweaker_c::TweakError_e eError;
eError = pcTweaker->TweakValue( "Shadowmap scale", 20.0F );
```

因此，调节变量只需一行代码，另外还需要几行处理错误的代码。取得被存储的的方式与此类似：

```
float vShadowmapScale;
eError = pcTweaker->GetValue( "Shadowmap scale", &vShadowmapScale );
```

### 1.18.4 图形用户界面

GUI 通常随项目而异，因此本文不打算对这一主题进行讨论。不过我将介绍一个已有的实现，让读者对该主题有一定的认识。在 **Equinox**（Funcom 的下一代 3D 引擎）中，我们实现了一个目录结构（如图 1.18.2 所示），在控制台中，可将该目录显示在屏幕的最上面。

```

Tweaker: Application
[...]
Caustics depth
Fog density
Fog end
Fog start
Linear fog
Physical water
Show Equinox logo
Show caustics
Show fog
Show landscape
Show sky
Show water
Table fog

```

图 1.18.2 GUI 的屏幕快照。用户可以在目录  
(在代码中为类别) 中上下移动, 并选择要调节的值

为调节值, 我们定义了可被指派给可调节变量的输入方法。这样, 我们可以创建专用的输入方法, 如角度调节器 (如图 1.18.3 所示)。

```

Tweaker: Graphics
Tweaker instance name: GraphicsTestInstance
AngleTweak 1/2
Type: float
Value = 56.649902
Limited to range <45.000000, 120.000000>
step = 0.750000, use +/-/space to modify

```




图 1.18.3 该专用输入方法让用户能够直观地调节角度

在保存和装载方面, 除了二进制瞬态图外, 我们还可以直接将 `#define` 语句中的所有可调节变量保存到头文件中。由于在应用程序的生命周期内, 变量的实例数可能发生变化, 我们只将第一个实例保存到头文件中。这种特性让我们能够只在调试版本中加入要调节的变量, 然后在发行版本中包含该头文件, 将变量初始化为最后调节得到的值。下面是一个例子, 它演示了如何采用这种方法来处理 `ShadowmapScale` 变量:

landscape\_tweakables.h:

```

...
#define SHADOWMAP_SCALE 43.5
...

```

landscape.cpp:

```

#include "landscape_tweakables.h"
...
m_vShadowmapScale = SHADOWMAP_SCALE;
...

```

### 1.18.5 附注

---

可以使用 RTTI `typeid()` 来替换前面介绍的类型信息代码。使用我们的类型信息代码有利也有弊。

优点是：

- 类型信息占用的空间小，因为只有使用类型信息的类才需要；
- 可以在 `TypeID_c` 类中添加特定的信息，例如一种将类型或指针装载并存储到 GUI 控件中的方式。

缺点是：

- 必须对每种类型使用宏，而 RTTI 能自动提供类型信息。

### 1.18.6 致谢

---

这里要感谢 Robert Golias，他提供了宝贵的帮助和建议，并实现了 Equinox 调节器的 GUI。该 GUI 证明了可以提供这样一个如此简单的界面。





## 1.19 生成真正的随机数

---

Pete Isensee, Microsoft

pkisensee@msn.com

计算机游戏在掷骰子、洗牌、模拟自然环境、产生逼真的物理现象和实现安全的多玩家事务时大量地使用随机数。计算机擅长生成伪随机数，但不擅长生成真正的随机数。伪随机数看起来是随机的，但在算法上是基于前一个随机数计算得到的。真正的随机数不但看起来是随机的，而且是不可预测、不重复、不确定的，这种随机数并不是将前一个随机数作为输入来生成的。本文介绍一种在软件中生成真正的随机数的方法。

### 1.19.1 伪随机

---

伪随机数序列总会重复，只要使用相同的种子，便可准确地再现。在游戏场景中，这将导致问题。来看游戏中一种常见的情况：使用当前时钟周期数（从计算机启动后经过的时钟周期）初始化随机数生成器（RNG）。假设玩家每次开始该游戏时都重新启动其游戏控制台，游戏的随机程度是由种子决定的，因此随机性的量级将小到难以接受的程度。

现在，假设使用伪随机数生成器来创建一个密钥，用于对多玩家游戏传输的数据进行加密。对于所有的公开密钥加密系统，其核心都是生成一个不可预测的随机数。使用伪随机数将无法保证安全，因为在初始状态已知的情况下，伪随机数是完全可以预测的，即很容易被破解。加密系统最薄弱的环节常常是其密钥的生成技术[Kelsey98]。

### 1.19.2 真正随机

---

真正的随机数满足以下条件：是随机的、分布不均匀、不可预测、不重复。对安全性而言，不可预测性至关重要。黑客即使对使用的算法非常了解，也无法预测结果[Scheneier96]。

生成随机数的理想方式是使用随机的物理源，如放射性衰变或热噪音。这样的设备有很多，例如[Walker(a)]。然而，PC和视频游戏控制台通常无法访问这样的设备。在没有硬件源的情况下，RFC 1750[Eastlake94]推荐使用的技术是：从大量不相关的源中获得随机输入，并使用一个强大的混合函数将它们混合在一起。通过从很多不相关的源那里获得输入（每个源都有一定量级的随机性），并充分地将它们混合起来，可获得一个混乱程度非

常高的值——一个真正的随机数。

### 1.19.3 随机输入源

---

在很多 PC 和游戏控制台中，可用于随机输入有：

- 系统日期和时间；
- 系统启动后经过的时间（最小时间单位）；
- 用户名或 ID；
- 计算机名称或 ID；
- CPU 寄存器的状态；
- 系统线程和进程的状态；
- 栈中的内容；
- 鼠标和游戏杆的位置；
- 最后 N 次键击或控制器输入之间的时间间隔；
- 最后 N 次键击或控制器输入；
- 内存状态（已分配的字节数、可用的字节数等）；
- 硬盘驱动器的状态（可用的字节数、已被占用的字节数等）；
- 最后的 N 个系统消息；
- GUI 状态（窗口的位置等）；
- 最后 N 个网络分组之间的时间间隔；
- 最后 N 个网络分组中的数据；
- 存储在主存储器、视频存储器等的半随机地址中的数据；
- 硬件标识符：CPU ID、硬盘驱动器 ID、BIOS ID、网卡 ID、视频卡 ID 和声卡 ID。

在给定的系统中，其中一些源保持不变，如用户 ID 和硬件 ID。之所有要包含这些值，是由于它们随机器而异，因此对于生成用于传输网络数据的密钥很有用。有些源的变化很小，例如在两次读取操作之间，硬盘驱动器的状态和内存负载的变化非常小。然而，每种输入都提供了一定程度的随机性，将它们混合起来后，可获得非常高的随机性。

从输入源获得的混乱程度越高，输出的随机性将越强。将诸如鼠标位置、键击和网络分组等存储到循环队列中很有用，这样可以将整个队列用作输入源。

### 1.19.4 硬件源

---

有些游戏平台能够访问具有随机性的物理源，这些物理源是非常棒的输入源。这样的物理源包括：

- 来自无源声卡（如麦克风插口）的输入；
- 来自摄像机的输入；
- 磁盘驱动器（硬盘驱动器、光驱、DVD）的寻道时间；
- Intel 810 芯片组的硬件 RNG（一个基于硅中热噪音的 RNG 实现）。

### 1.19.5 混合函数

生成真正的随机数时，强大的混合函数指的是这样的函数，以非线性方式将各项输入组合起来，从而生成输出。在输入的一个比特发生变化时，优秀的混合函数生成的输出中将有大约 50% 的比特也将变化。强大的混合函数有：

- DES（以及其他对称加密程序）；
- Diffie-Hellman（以及其他公共密钥加密程序）；
- MD5、SHA-1（以及其他哈希加密程序）。

安全的哈希函数（如 MD5）是理想的混合函数，原因在于：它们满足了优秀混合函数的基本要求、其安全性缺陷已被广泛地分析、速度通常比对称或非对称加密算法快、没有任何出口限制。公共密钥实现也被大量使用。

### 1.19.6 局限性

不同于生成伪随机数，在软件中生成真正的随机数的速度非常慢。为确保输出是真正随机的，必须对很多源进行取样。有些取样的速度非常慢，如从硬盘驱动器或声卡那里读取数据。另外，还必须使用复杂的算法对取得的输入进行混合。

相对于 PC，游戏控制台可选择的输入源更少，因此生成的结果的随机程度通常更低。然而，较新的控制台通常有某种硬盘驱动器（光驱、DVD、硬盘），它们可用作加密的硬件源。

结果的随机程度完全取决于输入样本的混乱程度。输入样本越多，每个样本的混乱程度越高，输出将越好。别忘了，这种算法被调用的频率越高，输出的随机性越低，这是因为输入中的比特变化更小。总而言之，这种技术并不能替代伪随机数生成器。这种技术适用于生成可持续使用几个小时或几天的随机数生成器种子值或网络会话密钥。

### 1.19.7 实现



附带光盘中有个真正的随机数生成器，它是用 C++ 编写的。这种算法的任何实现都不是独立于平台的。光盘中的实现只能用于 Win32 平台，但可以很容易地对其进行扩展，从而用于其他平台。该实现使用硬件源，如 Intel RNG 和声卡输入（如果它们可用的话）。出于简化和提高效率的目的，它并没有使用前面列出的所有输入，但使用的输入足以获得高度的随机性。

主要的功能是在 GenRand 对象中实现的，该对象位于名称空间 TrueRand 中。下面的范例演示了如何使用 GenRand 对象来生成真正随机的种子值：

```
#include "GenRand.h" // Genuine random number header
unsigned int nSeed = TrueRand::GenRand().GetRandInt();
```

下面的范例演示了如何生成一个会话密钥，以便进行安全的网络通信。Buffer 对象是一个简单的包装器（wrapper），它封装了 std::basic\_string<unsigned char>，提供了预留空间、追

加数据和记录样本缓冲区大小等所需的功能。

```
TrueRand::GenRand randGen;
TrueRand::Buffer bufSessionKey = randGen.GetRand();
```

函数 `GetRand()` 是该程序的核心，它获取随机输入样本，然后使用一个强大的混合函数来生成输出。该实现使用 MD5 哈希算法，因此得到的缓冲区长度为 MD5 哈希长度（16 Byte）。`mCrypto` 对象是一个封装了 Win32 Crypto API 的包装器，它提供了 MD5 哈希算法。

```
Buffer GenRand::GetRand()
{
    // Build sample buffer
    Buffer randInputs = GetRandomInputs();

    // Mix well and serve
    return mCrypto.GetHash( CALG_MD5, randInputs );
}
```

函数 `GetRandomInputs()` 是一个输入取样器，它返回一个缓冲区，其中包含大约 10K 的采样数据。可以根据需要，对该函数进行修改，以包含更多或更少的输入。由于该函数的运行时间随系统（驱动器、声卡）而异，因此可将硬件延迟（即函数开始和结束时的当前时间）作为随机输入源。

```
Buffer GenRand::GetRandomInputs()
{
    // For speed, preallocate input buffer
    Buffer randIn;
    randIn.reserve( GetMaxRandInputSize() );

    GetCurrTime( randIn );      // append time to buffer
    GetStackState( randIn );    // stack state
    GetHardwareRng( randIn );   // hardware RNG, if avail
    GetPendingMsgs( randIn );   // pending Win32 msgs
    GetMemoryStatus( randIn );  // memory load
    GetCurrMousePos( randIn );  // mouse position

    // . . . etc.

    GetCurrTime( randIn );      // random hardware latency

    return randIn;
}
```

最后，下面是一个输入取样函数。它取得当前时间，然后将数据附加到缓冲区对象 `mRandInputs` 中。`QueryPerformanceCounter()` 是 Windows 中精度最高的计时器，因此它提供了最大程度的随机性。在这里（及其其他很多情形下），我们可对 API 失败的情况置之不理，因为即使该函数失败，我们也将附加 `PerfCounter` 中的随机栈数据。

```
void GenRand::GetCurrTime( Buffer& randIn )
{
```

```
LARGE_INTEGER PerfCounter;  
QueryPerformanceCounter( &PerfCounter ); // Win32 API  
Append( randIn, PerfCounter );  
}
```

### 1.19.8 GenRand 的随机程度



有很多测试方法可以检查随机数的质量，其中之一是使用附带光盘中的 ENT 程序[Walker(b)]，该程序可对任何数据流（stream）进行一整套的测试。对 GenRand() 进行测试（不使用包括硬盘驱动器的寻道时间在内的任何硬件输入源，并使用 GetRandInt()生成 25 000 个随机整数）的结果如下：

- 混乱程度=7.998199 bit/byte;
- 使用最适合的压缩方法时，这个 100 000 byte 的文件不会减小；
- 100 000 个样本的 X 平方分布为 250.13，随机情况下，这个值应增加 50%；
- 数据字节的算术平均值为 127.4918（随机情况下，为 127.5）；
- $\pi$  的蒙特卡罗值为 3.15726293（误差为 0.5%）；
- 序列相关系数为 0.000272（完全不相关时为 0.0）。

上述结果表明，输出的随机性很高。例如，X 平方测试——最常用的随机性测试[KNUTH98]表明，该生成器的随机程度非常高。

### 1.19.9 参考文献

[Callas96] Callas, Jon, "Using and Creating Cryptographic-Quality Random Numbers," available online at [www.merrymeet.com/jon/usingrandom.html](http://www.merrymeet.com/jon/usingrandom.html), June 1996.

[Eastlake94] Eastlake, D., Network Working Group, et al, "Randomness Recommendations for Security," RFC 1750, available online at [www.faqs.org/rfcs/rfc1750.html](http://www.faqs.org/rfcs/rfc1750.html), December 1994.

[Kelsey98] Kelsey, J., et al, "Cryptanalytic Attacks on Pseudorandom Number Generators," available online at [www.counterpane.com/pseudorandom\\_number.html](http://www.counterpane.com/pseudorandom_number.html), March 1998.

[Intel99] Intel Corporation, "Intel Random Number Generator," available online at <http://developer.intel.com/design/security/rng/rng.htm>, 1999.

[Knuth98] Knuth, Donald, *The Art of Computer Programming, Volume 2: Seminumerical Algorithms, Third Edition*. Addison-Wesley. 1998.

[Schneier96] Schneier, Bruce, *Applied Cryptography, Second Edition*. John Wiley & Sons. 1996.

[Walker(a)] Walker, John, "HotBits: Genuine Random Numbers Generated by Radioactive Decay," available online at [www.fourmilab.ch/hotbits/](http://www.fourmilab.ch/hotbits/).

[Walker(b)] Walker, John, "ENT: A Pseudorandom Number Sequence Test Program," available online at [www.fourmilab.ch/random/](http://www.fourmilab.ch/random/).



## 1.20 使用 Bloom 过滤器来提高计算性能

Mark Fischer, Beach Software  
beach@beachsoftware.com

**假**设要在一个比特数组 (bit array) 中存储布尔信息——一个简单的假设。只需赋予比特数组中每个元素以特定含义, 然后赋给它一个值。在这种情况下, 存储 1 比特的信息需要一个比特。比特数组表示其相对值的准确度为 100%。当然, 当被存储的数据是面向数组的时 (如时空的瞬态), 这很好。然而, 如果数据不是面向线性瞬态的, 情况又如何呢?

### 1.20.1 Bloom 的方式

1970 年, Burton H. Bloom 在 *Communication of the ACM* 上发表了一种简单而灵巧的算法[Bloom70]。在这篇文章中, Bloom 建议使用一种“有可接受误差的哈希编码”算法来帮助字处理器检查文档进行大写和断字处理。与传统的一对一映射算法相比, 这种算法使用的空间更小、速度更快。使用这种算法, 只需使用一个简单的规则, 便可以检查大部分单词 (例如, 90%), 同时使用一个例外清单 (exception list) 来解决余下的一小部分单词, 该清单用来捕捉该算法指出某个单词可以被解决, 而实际上不可解决的情况。Bloom 这样做旨在缩短从缓慢的存储设备中查找数据所需的时间。

### 1.20.2 可能的情形

Bloom 过滤器 (Filter) 可以缩短完成相对昂贵且经常进行的计算所需的时间, 这是通过存储前一次计算中的一个布尔值来实现的。请看下面的情形, 我们需要改进它们的性能。

- 确定一个多边形是否在一个八叉树节点中可见;
- 确定对象是否在发生碰撞;
- 确定光线是否会射到对象上。

它们都属于一种更通用的情形。它们都涉及到昂贵的计算 (CPU、网络或其他资源), 而结果是一个布尔值 (通常为 false)。注意每种情形下可能使用的措词 (word) 至关重要, 因为如果 Bloom 过滤器测试返回 false (miss), 则 Bloom 过滤器将 100% 准确; 而仅当 Bloom 过滤器测试返回 true 时, 才可能不准确。

Bloom 过滤器可存储任何函数的真值结果。通常, 函数的参数是一个

指向字节数组的指针。如果要存储使用多个参数的函数返回的结果，可以将这些参数拼接成一个参数。当要求 100% 准确时，如果 Bloom 过滤器测试返回 true，必须计算最初的昂贵的函数，以确定该函数的绝对结果。

### 1.20.3 工作原理

Bloom 过滤器中有两个主要的函数，一个用于存储昂贵函数返回的布尔真值，另一个用于测试以前存储的布尔真值。存储函数将接受任何形式的输入，并相应地修改 Bloom 过滤器数组。测试函数接受的输入格式与存储函数相同，并返回一个布尔值。如果测试函数返回 false，则说明输入没有被存储函数存储过；如果测试函数返回 true，则说明输入以前可能被存储函数存储过。测试可能返回假阳性 (false positive)。如果要求 100% 准确，则必须执行原来的昂贵函数，以确定绝对值。可以添加一个传统的 Bloom 过滤器，这样可以存储昂贵函数返回的其他布尔真值，而不会删除以前存储的值。

### 1.20.4 定义

要获得高质量的 Bloom 过滤器，必须使用高质量的哈希函数，后者有时被称作消息摘要 (digest) 算法。任何高质量的哈希函数都管用，但笔者建议使用 RSA Security 公司的 MD5 消息摘要算法[RSA01]。该算法的源代码可从网上下载，另外 RFC 1321 中有有关该算法的文档。MD5 哈希函数使用一个字节数组中的 N 个字节，并生成一个 16 字节 (128 bit) 的返回值。该返回值是输入的散列 (hash)，这意味着输入中的任何位发生变化时，返回值都将发生巨大的变化。Bloom 将哈希函数的返回值称作 Bloom 过滤器键值 (key)。

Bloom 过滤器索引是通过将 Bloom 过滤器键值划分为指定长度的块而获得的。如果 Bloom 过滤器索引的长度为 16 bit，则一个 128 bit 的 Bloom 过滤器键值可被划分为 8 个完整的段。将键值划分为完整的段后，剩余的位将被丢弃。

Bloom 过滤器使用的相位 (phase) 数等于 Bloom 过滤器索引数，后者用于存储昂贵函数返回的布尔值。在这个例子中，使用了 128 bit 键值 (Bloom 过滤器索引的长度为 16 bit) 中的 3 个相位，余下的 5 个索引被丢弃。

Bloom 过滤器数组用于存储昂贵函数的布尔值。例如，如果 Bloom 过滤器索引的长度为 16 bit，则 Bloom 过滤器数组的长度将为  $2^{16}$  bit，即 64Kb (8KB)。数组越大，Bloom 过滤器测试的准确度越高。

Bloom 过滤器数组的 Bloom 过滤器饱和度 (Saturation) 指的是该比特数组中被设置为 true 的比特所占的百分比。当饱和度为 50% (即一半的比特被设置为 true) 时，Bloom 过滤器数组处于最佳状态。

### 1.20.5 范例 1

我们将使用 3 个相位 (索引长度为 16 bit) 来将函数参数 ("Mikano is in the park") 存储到一个 64Kb (8KB) 的数组中。在这个例子中，如果 Mikano 确实在公园里，则昂贵函数返



回 true。虽然这里使用的是一个字符串变量，但也可以使用任何变量格式。存储的昂贵函数参数的格式不会影响 Bloom 过滤器的性能、准确性以及使用的内存量。

首先，针对昂贵函数的参数执行哈希函数。我们假设该哈希函数返回一个 128 bit 的值：0x1002AB30001BF7877AB34D976A09667，则前 3 个 16 bit 的索引为 0x1002、0x7AB3、0x0001，而余下的索引将被丢弃。

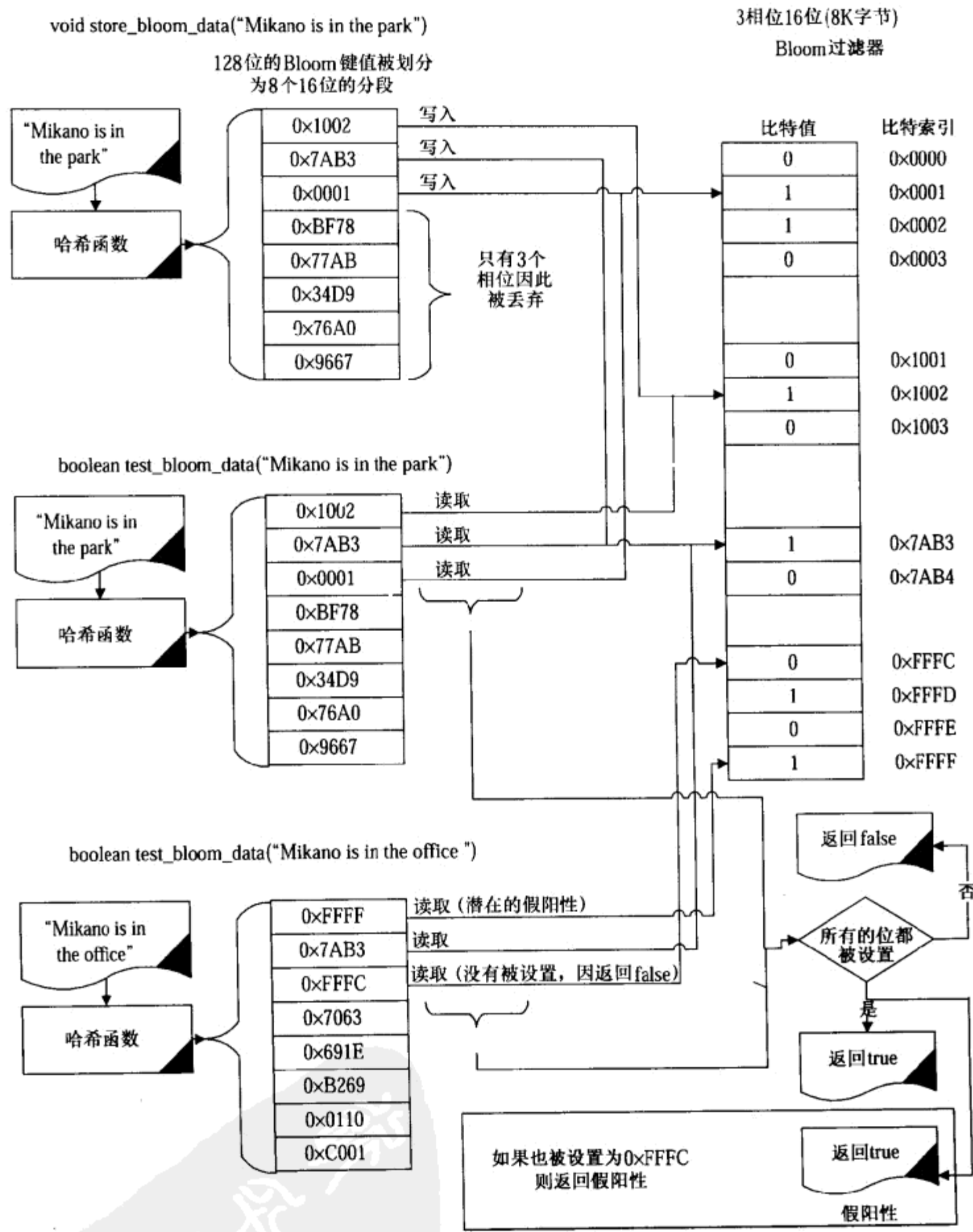


图 1.20.1 Bloom 过滤器的流程

```

// returns a pointer to 16 bytes of data that represent the hash
void * compute_hash( pData, nDataLength );

// returns the integer value for the bits at nIndex for nBitLength long
int get_index_value( void * pData, int nIndex, int nBitLength );

// tests a bit in the Bloom Filter Array.
// returns true if set otherwise returns false
boolean is_bit_index_set( int nIndexValue );

// sets a bit in the Bloom Filter Array
void set_bit_index( int nIndexValue );

void store_bloom_data( void * pData, int nDataLength )
{
    void *pHash;
    int nPhases = 3, nPhaseIndex = 0, nBitIndexLength = 16;

    // returns pointer to 16 bytes of memory
    pHash = compute_hash( pData, nDataLength );

    // now set each bit
    while ( nPhaseIndex < m_nPhases )
    {
        nIndexValue = get_index_value( pHash, nPhaseIndex, nBitIndexLength );

        // if bit is not set, we have a miss so return false
        set_bit_index( nIndexValue );

        nPhaseIndex++;
    }
}

boolean test_bloom_data( void * pData, int nDataLength )
{
    void *pHash;
    int nPhases = 3, nPhaseIndex = 0, nBitIndexLength = 16;

    // returns pointer to 16 bytes of memory
    pHash = compute_hash( pData, nDataLength );

    // now test each bit
    while ( nPhaseIndex < m_nPhases )
    {
        nIndexValue = get_index_value( pHash, nPhaseIndex, nBitIndexLength );

        // if bit is not set, we have a miss so return false
        if ( !is_bit_index_set( nIndexValue ) ) return( false );

        nPhaseIndex++;
    }

    // all bits are set so we have a probably hit.
    return( true );
}

```

从理论上说，不同的输入参数也可能返回同样的值，但这是不可能的。无论是哪种情况，该算法都管用。

使用相同的输入参数调用时，compute\_hash 总是返回16字节相同的数据

发现 false 位后，立刻返回 false。这说明该昂贵函数参数以前没有被存储过

图 1.20.2 Bloom 过滤器的基本用法

使用数据填充比特数组之前，重置 Bloom 过滤器数组（全部设置为 false）。然后，对于每一个索引，我们都将 Bloom 过滤器数组中相应的位设置为 true，而不管其以前的值是多少。填充该数组后，我们有时候需要将已经被设置为 true 的位设置为 true。这是导致测试 Bloom 过滤器数组时，结果为假阳性的原因所在（参见图 1.20.1）。

当需要检查 Bloom 过滤器数组，以确定以前是否存储了昂贵函数的参数时，步骤几乎与存储时相同，只不过是读取 Bloom 过滤器数组中的位，而不是将值写入到这些位中。如果被读取的任何位为 false，则说明该昂贵函数参数以前没有被存储到 Bloom 过滤器数组中；如果所有的位都为 true，则说明可能以前已经将该昂贵函数参数存储到数组中。在这种情况下，需要执行原来的昂贵函数，准确地确定布尔值（参见图 1.20.2）。

调整 Bloom 过滤器

为调整 Bloom 过滤器，需要确定相位数和索引的长度。通过修改这两个变量，可以改变 Bloom 过滤器的准确度和容量。一般来说，比特数组（N）越大，相位数越多，报告假阳性的可能性越小。Bloom 指出，当比特数组的饱和度为 50% 时，该算法的性能最高。统计数据表明，发生假阳性的概率为数组饱和度的相位数次幂。有一些可用于调整 Bloom 过滤器算法的公式。

计算假阳性概率的公式如下：

$$\text{percent\_false\_positive} = \text{saturation}^{\text{number\_of\_phases}}$$

上式也可表示为：

$$\text{number\_of\_phases} = \text{Log}_{\text{saturation}}(\text{percent\_false\_positive})$$

如果 Bloom 过滤器数组的容量处于最佳状态，即饱和度为 50%，则根据上述公式，可得到如表 1.20.1 的结果。

表 1.20.1 假阳性概率随相位数的变化情况

percent_false_positive	number_of_phases
50.00%	1
25.00%	2
12.50%	3
6.13%	4
3.13%	5
1.56%	6
0.78%	7
0.39%	8

例如，如果要求假阳性的概率低于 0.5%，则必须使用 8 个相位。在这种情况下，假阳性的最大概率为 0.39%。

接下来，我们来计算 Bloom 过滤器数组的长度：

```
array_bit_size = ( number_of_phases * max_stored_inputs ) / -ln(0.5)
```

array\_bit\_size 通常被上舍入（只进不舍）为与之最接近的 2 的整数次幂，即：

```
array_bit_size = 2n
```

最后，根据数组的长度来计算索引的长度：

```
array_bit_size = 2index_bit_size
```

1.20.6 范例 2

假设要存储最长为 9000 位的昂贵函数参数，且要求 Bloom 过滤器数组测试的返回值为 true 时，准确率为 95%。根据表 1.20.1 可知，要使准确率不低于 95%（即假阳性的概率不超过 5%），需要 5 个相位。

$5 \text{ phases} * 9000 \text{ expensive function parameters} / -\ln(0.5) = 64\,921 \text{ bit}$

将上述值上舍入为最接近的 2 的整数次幂，得到 64K bit（8KB）。由于  $2^{16} = 64\text{K}$ ，因此 index\_bit\_size 为 16 bit。

1.20.7 最后的说明

一种提高性能的方式是，像 Bloom 在其算法中所做的那样，使用例外清单来避免执行昂贵函数。例外清单中包含 Bloom 过滤器测试可能返回的所有假阳性情况。这可通过存储参数来计算，也可以在需要检测假阳性时动态地进行计算（如图 1.20.3 所示）。

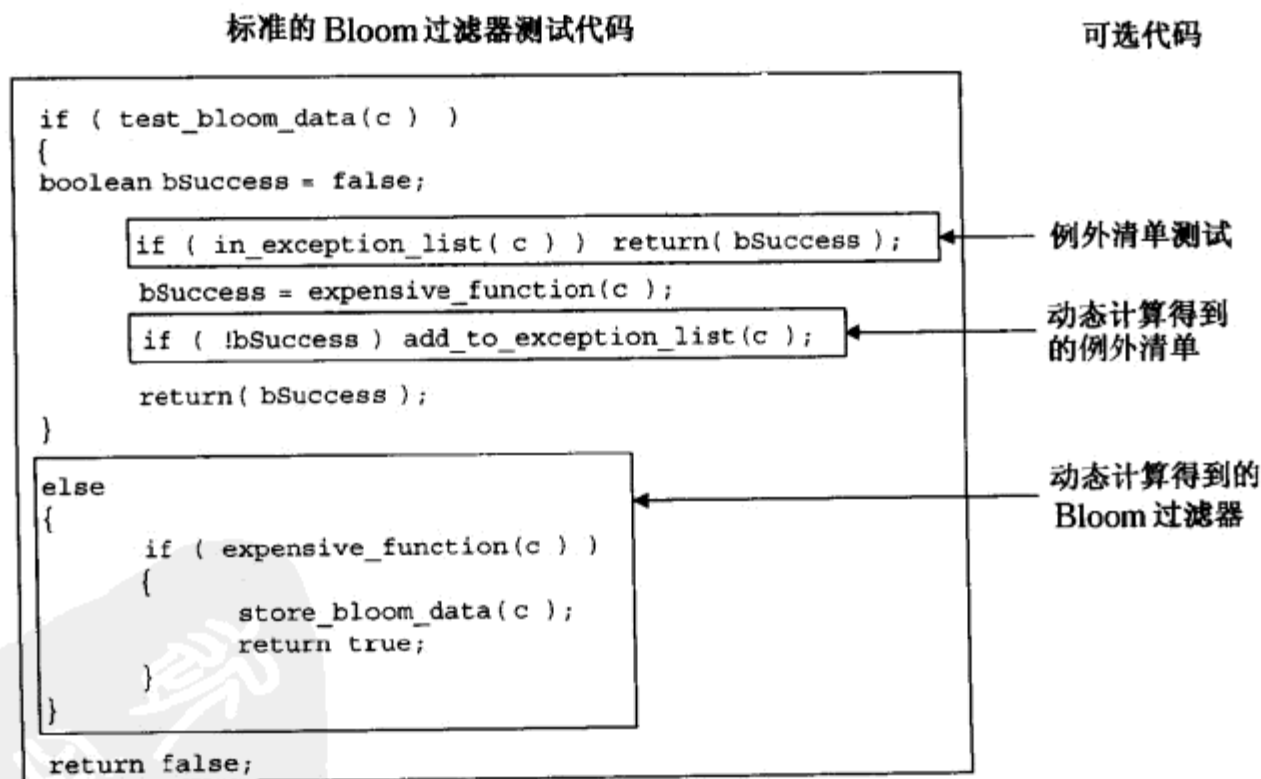


图 1.20.3 Bloom 过滤器的配置

另一种提高性能的方式是，动态地构建一个 Bloom 过滤器数组。如果昂贵函数参数的范围非常大，可以动态地计算 Bloom 过滤器，并针对每次的比特数组测试进行优化。通过动态地构建 Bloom 过滤器数组，常被测试的昂贵函数参数将被计算一次，而未被测试的函数参数不会浪费比特数组中的空间。

下面是 Bloom 过滤器的一些有趣的特征：

- 可通过按位 OR 运算，将两个 Bloom 过滤器数组合并起来；
- 不同的客户可共享同一个 Bloom 过滤器数组；
- 最优的 Bloom 过滤器数组是不可压缩的；
- 未满足的数组的可压缩性非常高；
- 通过将未知的位设置为 true，可以修复数组中的内存损坏（corruption）。

### 1.20.8 结论

---

Bloom 函数提供了一种方法，可以提高被频繁执行的昂贵函数的性能，但这是以占用更多的内存为代价的。虽然这种方法面世了很长时间，但仍是一种不常被使用的技术，虽然也存在例外情况，例如 Duane Wessels 开发的 Web 缓存程序 Squid ([www.squidcache.org/](http://www.squidcache.org/)) 大量使用了 Bloom 过滤器。在程序中添加 Bloom 过滤器算法所需的代码通常不超过 20KB。与大多数性能改进技巧一样，在优化阶段（主要功能实现后），将 Bloom 过滤器加入到项目中是个不错的主意。

### 1.20.9 参考文献

---

[Beach01] Beach Software, "Bloom Filters," available online at <http://beachsoftware.com/bloom/>, May 10, 2000.

[RSA01] RSA Security, "What Are MD2, MD4, and MD5," available online at [www.rsasecurity.com/rsalabs/faq/3-6-6.html](http://www.rsasecurity.com/rsalabs/faq/3-6-6.html), March 4, 2001.

[Flipcode01] Flipcode, "Coding Bloom Filters," available online at [www.flipcode.com/tutorials/tus\\_bloomfilter.shtml](http://www.flipcode.com/tutorials/tus_bloomfilter.shtml), September 11, 2000.

[Bloom70] Bloom, Burton H., "Space/Time Trade-Offs in Hash Coding with Allowable Errors," *Communications of the ACM*, Vol. 13, No.7 (ACM July 1970): pp. 422~426.



## 1.21 3Ds MAX 中的 Skin 导出器和动画工具包

Marco Tombesi

tombesi@infinito.it

在现代电影中，我们看到过一些奇妙的特效，其中的怪物（如恐龙）能够平稳地运动。我们知道这是如何实现的（使用诸如 3Ds MAX、LightWave、Maya 等），但我们如何将这种技术用于游戏中呢？



本文将介绍用于实现上述目的的一整套工具。我们将使用 3Ds MAX（和 Character Studio）中创建的动画人物，使之在游戏的实时场景中平稳地移动。在此期间，它将通过导出插件，并以自定义的数据格式被存储。本文只深入介绍导出方面，对于其他方面，附带光盘中的代码做出了充分的解释。

下面来介绍必须的步骤：

（1）动画是在 3Ds MAX 3.1（以下简称 MAX）和 Character Studio 2.2 中，使用 Biped 和/或 bones 以及体形修改器（Physique modifier）创建的。需要指出的是，虽然这些工具的新版本已经面世，但在任何新版本中，所需的算法都将与此类似。

（2）导出插件创建一种自定义的文件格式（.MRC），其中包括：

- 网格信息（顶点、法线）；
- 骨架结构（骨头树）；
- 每块骨头对网格顶点的影响值（influence value，权重），每个顶点可能受多块骨头的影响；
- 骨头（bone）动画：对于每块骨头，这可能由一组变换和旋转关键帧（key）组成（使用四元数），其中包括从动画开始到需要进行变换之间的准确时间（单位为 ms）。

（3）要读取.MRC 文件，需要一个可重用的 DLL，其中包含完整的源代码。

（4）渲染器（renderer）在关键帧之间插入帧（采用线性方式或更佳的方式），计算用于每块骨头的当前变换矩阵。这是通过使用从动画开始所过去的时间来实现的，从而实现了平稳的、独立于硬件的动画。

（5）对于每一帧，渲染器都重新计算每个顶点及其法线的位置。这种计算是根据当前变换矩阵以及每块骨头对特定顶点的影响值进行的。大多数矩阵运算都可使用图形硬件的变换和光照（lighting）特性来完成（例如 GeForce 和 Radeon 卡）。





在 MAX 中，并没有对如何使用插件导出动画数据做详细的说明。虽然网上有很多文章对蒙皮 (skinning) 技术进行了介绍，但很少实际解决了导出数据的问题。请阅读并研究附带光盘中有关本文的项目目录中的源代码和 Readme.txt 文件。更详细的信息，请参阅作者的网页[Tombesi01]，该网页将针对 MAX 4 进行更新。

本文基于一个层次式骨骼结构 (使用 Character Studio 2.2 创建的骨头树或 Biped) 创建一个小型的多边形网格 (mesh) —— 大约包含 5000 个三角形。在 MAX 中，该网格是单个可被选中的对象。然后，使用 Physique，根据前面创建的 Biped，将该网格分解。人物动画应在 Biped 上创建。

### 1.21.1 导出

首先，我们需要一种文件格式规范。

#### 1. MRC 文件格式



这是一种简单的文件格式，用于实现本文的目标。它支持法线、骨头 (bone)、顶点、权重 (weight) 和动画关键帧。图 1.21.1 对其进行了说明，有关技术细节，请参阅附带光盘中的代码。

#### 2. 使用 MAX SDK 导出为 MRC

如果您不熟悉插件开发，也不了解 MAX 的工作原理，请参阅 MAX SDK 文档。具体地说，请阅读以下部分：



- DLL, Library Functions, and Class Descriptors;
- Fundamental Concepts of the MAX SDK;
- Must Read Sections for All Developers;
- Nodes;
- Geometry Pipeline System;
- Matrix Representations of 3D Transformations.

#### 3. 使用节点

在我们的导出插件中，必须从 SceneExport 派生出一个类，并实现一些虚拟函数，其中之一是主导出例程。

```
class MRCexport : public SceneExport {
public:
    // Number of extensions supported
    int ExtCount() {return 1;}
}
```



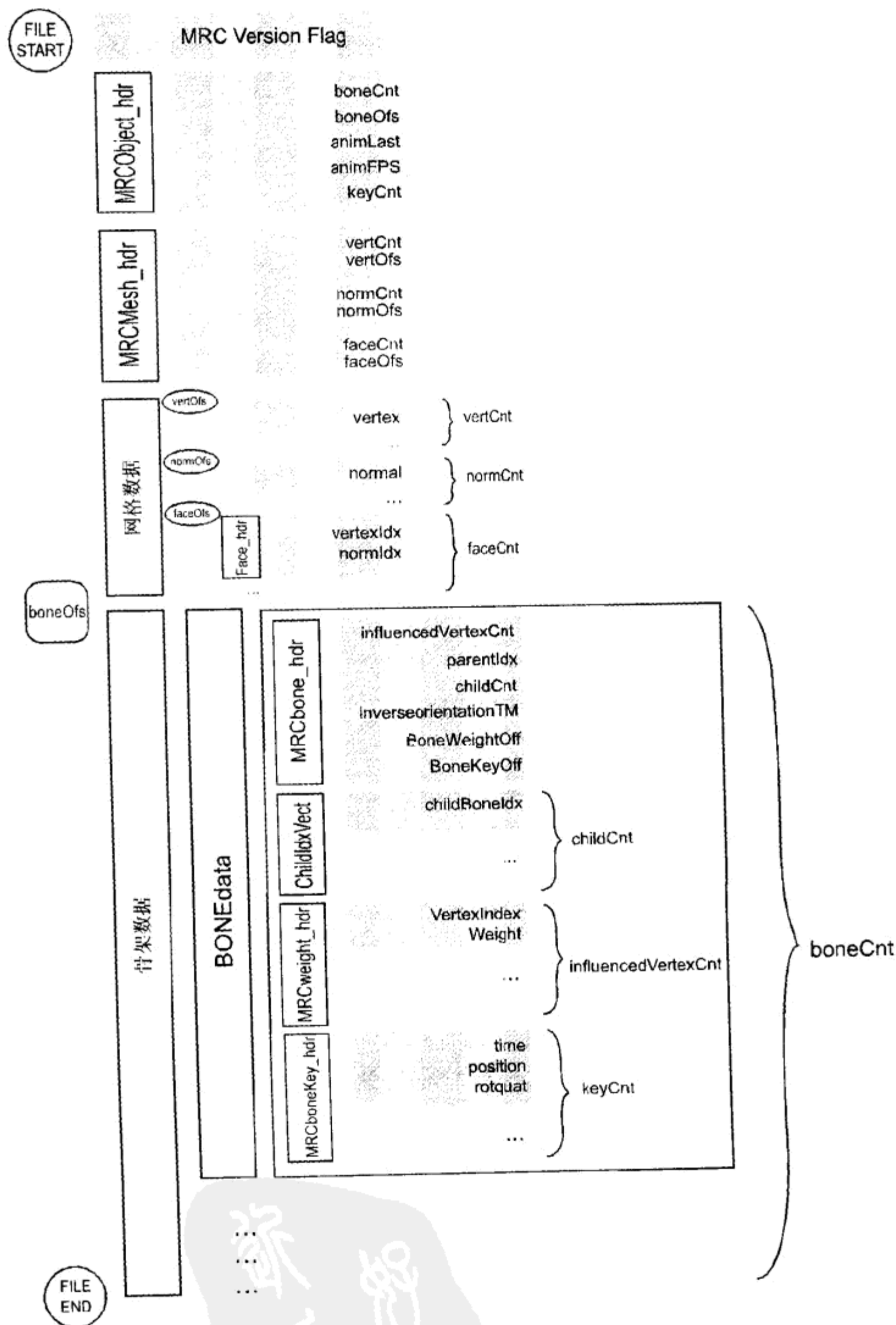


图 1.21.1 MRC 文件格式

```

// Extension("MRC")
const TCHAR * Ext(int n) {return _T("MRC");}
...

// Export to an MRC file
int DoExport(const TCHAR *name,
             ExpInterface *ei,
             Interface *i,
             BOOL suppressPrompts=FALSE,
             DWORD options=0);

//Constructor/Destructor
MRCexport();
~MRCexport();
...
};

```

要存取场景 (scene) 数据, 需要一个从 MAX 到主导出例程的接口 (插件的入口)。对于 MAX 中的每个对象, 全局场景图中都有一个节点, 而每个节点都有父节点 (RootNode 除外), 还可能有一些子节点。我们可以访问根节点, 然后遍历整个层次结构; 如果导出前, 用户选中了某个节点, 则可以直接访问它。

```

INode* pNode = i->GetSelNode(0);
INode* const pRoot = i->GetRootNode();

```

为在节点结构中导航, 我们编写了下述方法:

```

Int count = pNode->NumberOfChildren();
INode* pChNode = pNode->GetChildNode(i);

```

节点可以表示任何东西, 因此我们需要通过节点的类标识符 (Class\_ID 或 SuperClassID) 来确定其对象类型, 然后对其做相应的强制类型转换。就我们的目标而言, 需要检查节点是几何体对象 (网格) 还是骨头 (Biped 或 bone)。

```

bool IsMesh(INode *pNode)
{
    if(pNode == NULL) return false;
    ObjectState os = pNode->EvalWorldState(0);
    if(os.obj->SuperClassID() == GEOMOBJECT_CLASS_ID)
        return true;
    return false;
}

bool IsBone(INode *pNode)
{
    if(pNode == NULL) return false;
    ObjectState os = pNode->EvalWorldState(0);
    if (!os.obj) return false;

    if(os.obj->ClassID() == Class_ID(BONE_CLASS_ID, 0))

```

```

        return true;

    if(os.obj->ClassID() == Class_ID(DUMMY_CLASS_ID, 0))
        return false;

    Control *cont = pNode->GetTMController();
    //other Biped parts
    if( cont->ClassID() == BIPSLAVE_CONTROL_CLASS_ID ||
        //Biped root "Bip01"
        cont->ClassID() == BIPBODY_CONTROL_CLASS_ID
    ) return true;
    return false;
}

```

上述范例说明了如何在 MAX 节点中导航，以及如何检查它们表示的是什么。获得网格节点后，需要得到所需的顶点数据。

#### 4. 获取网格数据

为方便后面的操作，我们将以全局坐标空间的方式存储所有的顶点数据。MAX 对象的坐标位于对象空间中，因此我们需要一个变换矩阵，应用于每个顶点和网格的法线。

在动画期间，我们可以使用 `GetObjectTM(TimeValue time)` 来获得全局变换矩阵。这个矩阵用于将顶点从对象空间变换到全局空间，也可用来获得网格顶点的全局空间坐标。为此，可以将顶点在对象空间中的坐标与该方法返回的矩阵相乘（在 MAX 中，为自右乘）。我们感兴趣的是动画开始时的网格数据，因此 `TimeValue` 为 0。

```
Matrix3 tm = pNode->GetObjectTM(0)
```



MAX 使用行向量  $1 \times 3$  和  $4 \times 3$  矩阵，因此要对向量进行变换，必须将其同矩阵进行自左乘。

向量和其他数据并非静态存储，而是每次动态地进行计算。要存取数据，必须首先计算几何管道（*geometry pipeline*），并指定要获得什么时候的对象状态。

MAX 有一个修改器栈系统，其中的每个对象都是修改链的结果。我们以参数化的基本几何体（如立方体）为基础，对其依次应用栈中的修改器，从而得到最终的物体。这就是物体管道，我们将使用其结果。最后得到的对象是一个 `DerivedObject`，它包含用于在修改器栈中进行导航的方法。

要获取某一动画时刻的结果，必须首先取得一个 `ObjectState`，这是通过调用方法 `EvalWorldState` 实现的。这将使 MAX 应用管道中的每个修改器。

```
ObjectState os = pNode->EvalWorldState(0);
```

`ObjectState` 包含一个指针，它指向管道中的对象，获得该对象后，便可以获得网格数据。为此，我们必须将该通用对象转换为一个几何对象，后者有一个用于构建网格表示的方法。

```
Mesh& mesh = *(((GeomObject*)os.obj)->GetRenderMesh(0, pNode, ...));
```

现在，存取 mesh 成员并将顶点、面和法线放到内存中将很容易。然后，就可以将它们写入到文件中了。为此，需要使用下述方法：Mesh::getNumVerts()、Mesh::getNumFaces()、Mesh::getVert(i)和 Mesh::getNormal(i)。

程序清单 1.21.1 演示了如何将网格数据导出到文件中。

## 5. 获取骨架结构

现在，我们需要将骨架的层次结构写入到一个输出数据文件中。我们从根节点开始，按深度优先的方式遍历整个树，对于其中的每块骨头，都执行几项操作。首先，给它的每个直接子节点和父节点指定一个索引，然后获取骨头的定向矩阵（orientation matrix）。

```
tm = pNode->GetNodeTM(0);
tm.Invert();
```



上述矩阵并非对象矩阵，虽然很相似。该矩阵与节点的枢轴点（可能不是对象的原点）相关。更准确的信息，请参阅 SDK 文档。我们将使用该矩阵将每个网格顶点从全局空间变换到相应的骨头空间，这样它便可以随骨头一起移动。由于我们必须将顶点同该矩阵的逆矩阵相乘，因此现在可以求它的逆矩阵，并保存渲染时间。

## 6. 获取骨头的影响值

这是本文最激动人心的部分——获取分配给顶点的骨头及其影响值（权重）。当多块骨头影响同一个顶点时，权重至关重要；而网格将被如何分解，取决于这两个因素（有关其中的理论，请参阅[Woodland00]）。这种分配应使用 Character Studio 2.2 中的 Physique 修改器来完成。有关修改器接口的信息，请参阅 Character Studio 自带的头文件 Phyexp.h。

首先，我们必须在要导出的对象的节点上找到 Physique 修改器（这是前面用来获得网格顶点数据的节点）。为此，可以访问被引用的 DerivedObject，然后扫描栈中每个被应用的修改器，直到找到 Physique 修改器（通过检查 Class\_ID）。

```
Modifier* GetPhysiqueMod(INode *pNode)
{
    Object *pObj = pNode->GetObjectRef();
    if(!pObj) return NULL;

    // Is it a derived object?
    while(pObj->SuperClassID() == GEN_DERIVOB_CLASS_ID)
    {
        // Yes -> Cast
        IDerivedObject *pDerivedObj =
            static_cast<IDerivedObject*>(pObj);

        // Iterate over all entries of the modifier stack
        int ModStackIndex = 0;
        while(ModStackIndex < pDerivedObj->NumModifiers())
        {
```

```

        // Get current modifier
        Modifier* pMod = pDerivedObj->
            GetModifier(ModStackIndex);

        // Is this Physique?
        if(pMod->ClassID() ==
            Class_ID(PHYSIQUE_CLASS_ID_A,
                PHYSIQUE_CLASS_ID_B))
            return pMod;

        // Next modifier stack entry
        ModStackIndex++;

    }
    pObj = pDerivedObj->GetObjRef();
}
// Not found
return NULL;
}

```

现在，我们进入了骨头分配阶段（有关代码，请参阅程序清单 1.21.2）。有了 Physique 修改器后，便可以获得其接口（IPhysiqueExport）并访问对象的 Physique 场境（context）接口。该接口中包含我们需要的所有方法。受修改器影响的每个顶点都有一个 IPhyVertexExport 接口。要获得该接口，以访问它的方法，可以调用 Physique 场境接口的 GetVertexInterface(i) 方法。

我们必须检查顶点是受一块骨头还是多块骨头影响（分别为 RIGID\_TYPE 和 RIGID\_BLENDED\_TYPE）。对于前一种情况，权重为 1，因此只需找出一块骨头（调用顶点接口的 GetNode 方法）；对于后一种情况，我们必须找出分配给该顶点的所有骨头，并获取每块骨头的权重，这是通过调用顶点接口的 GetWeight(j) 方法实现的，其中 j 表示影响该顶点的第 j 块骨头。最后，我们必须释放所有的接口。

## 7. 获取骨头的动画关键帧（bone animation key）

这一步很简单。每过指定的时间间隔（默认为 100 ms），获取每个骨头的变换矩阵。在 MAX SDK 中，使用的时间单位为 tick（每秒为 4800 tick），因此我们必须对单位进行转换。然后，我们使用下面的方法：

```
tm = pNode->GetNodeTM(timeTicks);
```

与存储整个变换矩阵相比，只存储平移值（3 个浮点数）和旋转数据（4 个浮点数）的效率更高，因此，我们从变换矩阵中提取一个位置向量和一个单位四元素。

```
Point3 pos = tm.GetTrans();
Quat quat(tm);
```

将所有数据收集到内存中后，便可以将它们存储到磁盘中（文件格式为 MRC）。接下来介绍如何在游戏中使用该文件来实现平稳的动画。



## 8. 使用 MRC 文件：绘制循环

在我们的应用程序中，对于要显示的每一帧，都需要依次执行以下步骤。

## 9. 获得准确的时间

为确保动画平稳且独立于处理器，必须获取系统时间。我们将通过遍历骨头树来更新骨架结构，对于其中的每块骨头，将通过在两个样本关键帧之间进行线性插值来计算出当前的变换矩阵。为确定要在哪两个样本关键帧之间进行插值，必须知道从动画开始，到当前过去的动画时间（单位为 ms）。

## 10. 移动骨架

我们通过在选择的样本关键帧之间进行线性插值和四元数插值（SLERP）来确定骨头的位置和旋转（当前时间应该位于样本时间之间）。然后，根据平移和旋转数据创建当前骨头的动画矩阵。有关其中涉及的数学知识（尤其是四元数的计算），请参阅《游戏编程精粹 1》[Shankel00]。为充分利用图形硬件，我们将使用 OpenGL 函数来完成所有的矩阵计算。这样可以利用高级硬件特性，如变换和光照（lighting），因此性能将高得多。

## 11. 重新计算 Skin

移动骨架后，需要根据顶点的权重分配情况，相应地对网格进行分解。有关这一主题的简要介绍，请参阅[Woodland00]。可按骨头的重要性依次检查顶点，以深度优先的方式遍历骨头树，并对每块骨头执行下述操作。对于骨头影响的每个顶点，我们参照骨头的坐标系（乘以骨头定向矩阵的逆矩阵），然后通过骨头当前的动画矩阵对其进行变换。接下来，我们将顶点坐标与该骨头对它的影响值（权重）相乘，并将结果与存储在临时缓冲区中的顶点值相加。这样，缓冲区中便包含当前顶点在此时的动画中的坐标。最后，我们使用顶点数组（以获得更高的性能）绘制计算得到的网格。

### 程序清单 1.21.1 将网格数据导出到文件中

```
bool ExportMesh (INode* pNode, FILE *out)
{
    ...
    MRCmesh_hdr mHdr;
    Matrix3 tm = pNode->GetObjectTM(0);
    ObjectState os = pNode->EvalWorldState(0);
    int needDelete;
    Mesh& mesh = *((GeomObject*) os.obj)->GetRenderMesh (0, pNode, ...);
    ...

    // write the mesh vertices
    mHdr.vertCnt = mesh.getNumVerts();
    for(int i = 0; i < mHdr.vertCnt; i++)
    {
        Point3 pnt = mesh.getVert(i) * tm;    //premultiply in MAX
        ...
    }
}
```

```

}

// write vertex normals
mesh.buildNormals();
mHdr.normCnt = mesh.getNumVerts();
for(i = 0; i < mHdr.normCnt; i++)
{
    Point3 norm = Normalize(mesh.getNormal(i));
    ...
}

// build and write faces
mHdr.faceCnt = mesh.getNumFaces();
for(i = 0; i < mHdr.faceCnt; i++)
{
    MRCface_hdr fHdr;
    fHdr.vert[0] = mesh.faces[i].v[0];
    fHdr.vert[1] = mesh.faces[i].v[1];
    fHdr.vert[2] = mesh.faces[i].v[2];
    ...
}
...
}

```

### 程序清单 1.21.2 读取骨头分配信息

```

bool GetPhysiqueWeights(INode *pNode, INode *pRoot,
                        Modifier *pMod, BoneData_t *BD)
{
    // create a Physique Export Interface for given Physique Modifier
    IPhysiqueExport *phyInterface = (IPhysiqueExport*)
        pMod->GetInterface(I_PHYINTERFACE);
    if(phyInterface)
    {
        // create a ModContext Export Interface for the specific
        // node of the Physique Modifier
        IPhyContextExport *modContextInt = (IPhyContextExport*)
            phyInterface->GetContextInterface(pNode);

        // needed by vertex interface (only Rigid supported by now)
        modContextInt->ConvertToRigid(TRUE);

        // more than a single bone per vertex
        modContextInt->AllowBlending(TRUE);
        if(modContextInt)
        {
            int totalVtx = modContextInt->GetNumberVertices();
            for(int i = 0; i < totalVtx; i++)
            {
                IPhyVertexExport *vtxInterface = (IPhyVertexExport*)
                    modContextInt->GetVertexInterface(i);
            }
        }
    }
}

```

```

if(vtxInterface)
{
    int vtxType = vtxInterface->GetVertexType();
    if(vtxType == RIGID_TYPE)
    {
        INode *boneNode = ((IPhyRigidVertex*)vtxInterface)
            -> GetNode();
        int boneIdx = GetBoneIndex(pRoot, boneNode);
        Insert
        // Build vertex data
        MRCweight_hdr wdata;
        wdata.vertIdx = i;
        wdata.weight = 1.0f;

        //Insert into proper bonedata
        BD[boneIdx].weightsVect.push_back(wdata);

        // update vertexWeightCnt for that bone
        BD[boneIdx].Hdr.vertexCnt
            = BD[boneIdx].weightsVect.size();
    }
    else if(vtxType == RIGID_BLENDED_TYPE)
    {
        IPhyBlendedRigidVertex *vtxBlendedInt =
            (IPhyBlendedRigidVertex*)vtxInterface;
        for(int j = 0; j < vtxBlendedInt->GetNumberNodes()
            ;j++)
        {
            INode *boneNode = vtxBlendedInt->GetNode(j);
            int boneIdx = GetBoneIndex(pRoot, boneNode);

            // Build vertex data
            MRCweight_hdr wdata;
            wdata.vertIdx = i;
            wdata.weight = vtxBlendedInt->GetWeight(j);

            // check vertex existence for this bone
            bool notfound = true;
            for (int v=0; notfound
                && v < BD[boneIdx].weightsVect.size(); v++)
            {
                // update found vert weight data for that
                // bone
                if ( BD[boneIdx].weightsVect[v].vertIdx
                    == wdata.vertIdx )
                {
                    BD[boneIdx].weightsVect[v].weight
                        += wdata.weight;
                    notfound = false;
                }
            }
        }
    }
}

```

```

        if (notfound)
        {
            // Add a new vertex weight data into proper
            // bonedata
            BD[boneIdx].weightsVect.push_back(wdata);

            // update vertexWeightCnt for that bone
            BD[boneIdx].Hdr.vertexCnt
                = BD[boneIdx].weightsVect.size();
        }
    }
}

phyInterface->ReleaseContextInterface(modContextInt);
}

pMod->ReleaseInterface(I_PHYINTERFACE, phyInterface);
}

return true;
}

```

### 1.21.2 参考文献

- [Discreet00] Max SDK Plug-in development documentation: *SDK.HLP*
- [Tombesi01] Tombesi Marco's Web page: <http://digilander.iol.it/baggior/>
- [Woodland00] Woodland, Ryan, "Filling the Gaps—Advanced Animation Using Stitching and Skinning," *Game Programming Gems*. Charles River Media 2000; pp. 476~483.
- [Shankel00] Shankel, Jason, "Matrix-Quaternion Conversions" and "Interpolating Quaternions," *Game Programming Gems*. Charles River Media 2000; pp. 200~213.

## 1.22 在视频游戏中使用 Web 摄像机

Nathan d'Obrenan, Firetoad Software  
nathand@firetoads.com

**当**前的大部分游戏都允许多个玩家一起玩,但玩家之间的交互仅限于少量的文本消息。如果能够在晃过对手,并到达终点之前或对手被您发射的导弹炸死时,看到对手的表情,那该多好。Web 摄像机让您能够实现这种功能;同时随着高速 Internet 逐渐成为标准配置,将更多的数据发送给更多的客户是可行的。

本文演示一种在游戏中实现 Web 摄像机功能的简单方法。我们将使用 Video for Windows 来捕获 Web 摄像机的数据,因此需要使用 Windows,以执行初始化 Web 摄像机的函数。我们将介绍大量实现快速图像采集、运动检测的途径以及两个操纵图像的例程。最后的应用程序将具备完整的 Web 摄像机功能,其运行和交互的帧速率适度。

### 1.22.1 初始化 Web 摄像机捕获窗口



下面的代码演示了如何使用 Video for Windows 在应用程序中建立一个 Web 摄像机窗口。使用硬件厂商的视频驱动程序时,无需太多的错误检查和处理代码(有关实现的详细情况,请参阅附带光盘中的源代码)。

```
// Globals
HWND hWndCam = NULL;
BOOL cam_driver_on = FALSE;
int wco_cam_width = 160, wco_cam_height = 120;
int wco_cam_updates = 400, wco_cam_threshold = 120;

// WEBCAM_INIT
void webcam_init(HWND hWnd)
{
    // Set the window to be a pixel by a pixel large
    hWndCam = capCreateCaptureWindow(appname,
                                     WS_CHILD | WS_VISIBLE |
                                     WS_CLIPCHILDREN |
                                     WS_CLIPSIBLINGS,
                                     0, 0,
                                     1, 1,
                                     hWnd,
```

```
0);

if(hWndCam)
{
    // Connect the cam to the driver
    cam_driver_on = capDriverConnect(hWndCam, 1);

    // Get the capabilities of the capture driver
    if(cam_driver_on)
    {
        capDriverGetCaps(hWndCam, &caps, sizeof(caps));

        // Set the video stream callback function
        capSetCallbackOnFrame(hWndCam, webcam_callback);

        // Set the preview rate in milliseconds
        capPreviewRate(hWndCam, wco_cam_updates);

        // Disable preview mode
        capPreview(hWndCam, FALSE);

        // Initialize the bitmap info to the way we want
        capwnd.bmiHeader.biSize = sizeof(BITMAPINFOHEADER);
        capwnd.bmiHeader.biWidth = wco_cam_width;
        capwnd.bmiHeader.biHeight = wco_cam_height;
        capwnd.bmiHeader.biPlanes = 1;
        capwnd.bmiHeader.biBitCount = 24;
        capwnd.bmiHeader.biCompression = BI_RGB;
        capwnd.bmiHeader.biSizeImage = wco_cam_width*wco_cam_height*3;
        capwnd.bmiHeader.biXPelsPerMeter = 100;
        capwnd.bmiHeader.biYPelsPerMeter = 100;

        if(capSetVideoFormat(hWndCam, &capwnd,
                               sizeof(BITMAPINFO)) == FALSE)
        {
            capSetCallbackOnFrame(hWndCam, NULL);
            DestroyWindow(hWndCam);
            hWndCam = NULL;
            cam_driver_on = FALSE;
        }
        else
        { // Assign memory and variables
            webcam_set_vars();
            {
                glGenTextures(1, &webcam_tex.gl_bgr);
                glBindTexture(GL_TEXTURE_2D, webcam_tex.gl_bgr);
                glTexImage2D(GL_TEXTURE_2D, 0, 3, webcam_tex.size,
                             webcam_tex.size, 0, GL_BGR_EXT,
                             GL_UNSIGNED_BYTE, webcam_tex.bgr);

                glTexParameteri(GL_TEXTURE_2D, GL_TEXTURE_WRAP_S, GL_REPEAT);
            }
        }
    }
}
```



```

glTexParameteri(GL_TEXTURE_2D, GL_TEXTURE_WRAP_T, GL_REPEAT);

glTexParameterf(GL_TEXTURE_2D, GL_TEXTURE_MIN_FILTER,
                GL_LINEAR);
glTexParameterf(GL_TEXTURE_2D, GL_TEXTURE_MAG_FILTER,
                GL_LINEAR);
}

{
glGenTextures(1, &webcam_tex.gl_grey);
glBindTexture(GL_TEXTURE_2D, webcam_tex.gl_grey);
glTexImage2D(GL_TEXTURE_2D, 0, 1, webcam_tex.size,
             webcam_tex.size, 0, GL_LUMINANCE,
             GL_UNSIGNED_BYTE, webcam_tex.greyscale);

glTexParameteri(GL_TEXTURE_2D, GL_TEXTURE_WRAP_S, GL_REPEAT);
glTexParameteri(GL_TEXTURE_2D, GL_TEXTURE_WRAP_T, GL_REPEAT);

glTexParameterf(GL_TEXTURE_2D, GL_TEXTURE_MIN_FILTER,
                GL_LINEAR);
glTexParameterf(GL_TEXTURE_2D, GL_TEXTURE_MAG_FILTER,
                GL_LINEAR);
}
}
}
else
{
    cam_driver_on = FALSE;
}
}

```

上述函数获取 Web 摄像机窗口的句柄，我们通过函数 `capCreateCaptureWindow()` 来捕获该窗口。然后，我们初始化该窗口的属性（如大小），并确定是否使之可见。在我们的程序中，确实希望它可见，但我们将把它的大小设置为  $1 \times 1$ ，这样它几乎不可见。之所以这样做，是因为我们并不想显示该图像子窗口，但我们想通过回调函数来获取来自 Windows 的数据更新。

接下来，我们取得驱动程序的信息，设置回调函数（后面将更详细地介绍）和 Web 摄像机每秒的更新次数，然后重置所有的变量。然后，检查驱动程序是否能够返回我们更感兴趣的标准位图信息。如果能够，则初始化运动缓冲区的内存和 OpenGL 纹理（texture）。在决定该纹理应为多大时，使用了一个小窍门，这将在后面派上用场。我们根据 Web 摄像机窗口的高度，计算所需的内存量，并分配内存为：与计算机最接近，且大于它的 2 的整数次幂。虽然分配的内存量比存储 Web 摄像机图像所需的多，但这样可以避免调整纹理的大小，而这种操作是很昂贵的。这是通过使用 `memcpy()` 将内存复制到缓冲区实现的——代价是 Web 摄像机图像的精度有所降低。

## 1. 获取数据

初始化视频窗口后，我们需要获取 Web 摄像机每一帧的数据。为让 Windows 知道应将数据发送给哪个回调函数，我们必须调用 `capSetCallbackOnFrame()`，并将回调函数的地址传递给它。当 Windows 认为需要更新 Web 摄像机时，它将使用 `VIDEOHDR` 结构将位图信息传递给我们。

在我们的程序中，我们将让回调函数来处理所有的 Web 摄像机数据，以决定是否要使用这些数据来创建一个纹理。我们可以将所有这些数据传递给 `webcam_calc_movement()` 做进一步的处理，该函数将判断这些数据相对于最后一帧是否有足够大的变化，如果是，则更新纹理。

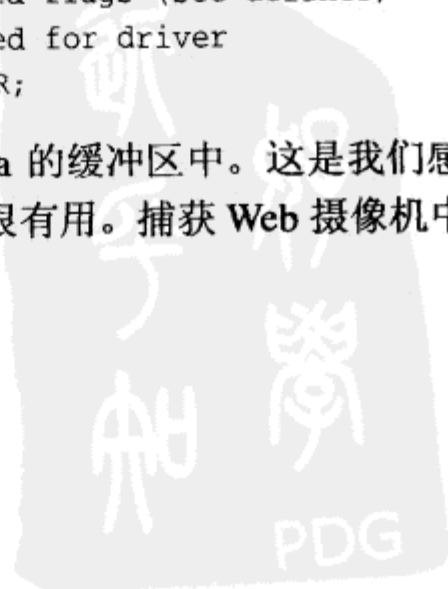
```
// WEBCAM_CALLBACK
// Process video callbacks here
LRESULT WINAPI webcam_callback(HWND hwnd, LPVIDEOHDR video_hdr)
{
    // Calculate movement based off of threshold
    if(webcam_calc_movement(video_hdr,
                            webcam_tex.delta_buffer,
                            wco_cam_width,
                            wco_cam_height,
                            webcam_tex.size,
                            wco_cam_threshold))
    {
        webcam_make_texture(video_hdr, wco_cam_rendering);
    }

    return TRUE;
}
```

Windows defines the `LPVIDEOHDR` structure as:

```
typedef struct videohdr_tag
{
    LPBYTE    lpData;           // pointer to locked data buffer
    DWORD     dwBufferLength;    // Length of data buffer
    DWORD     dwBytesUsed;       // Bytes actually used
    DWORD     dwTimeCaptured;   // Milliseconds from start of stream
    DWORD     dwUser;            // for client's use
    DWORD     dwFlags;           // assorted flags (see defines)
    DWORD     dwReserved[4];     // reserved for driver
} VIDEOHDR, NEAR *PVIDEOHDR, FAR *LPVIDEOHDR;
```

Windows 将 Web 摄像机数据保存在名为 `lpData` 的缓冲区中。这是我们感兴趣的主要变量，但 `dwTimeCaptured` 和一些标记 (flag) 也可能很有用。捕获 Web 摄像机中的数据后，我们对其进行检测，看它是否有用。



## 2. 运动检测

现在，需要清除变化不大的帧，以避免对纹理做不必要的更新。在诸如 OpenGL 等的 3D API 中，更新纹理的速度非常慢。

下述源代码对增量（delta）缓冲区进行比较，并根据是否超过指定的阈值而返回 true 或 false。注意，在超过域值的情况下，提前返回可进一步优化该函数；然而这将妨碍我们以后使用增量缓冲区。Ectosaver[Firetoad00]使用增量运动（delta movement）的这些无符号字节来计算它导致的浪幅（amplitude of the wave），并确定何时没有移动。

```
// GLOBALS
unsigned char wco_cam_threshold=128; // This is a good amount (0-255)

// WEBCAM_CALC_MOVEMENT
// This is a simple motion detection routine that determines if you've
// moved further than the set threshold
BOOL webcam_calc_movement(LPVIDEOHDR video_hdr,
                           unsigned char *delta_buff,
                           int webcam_width, int webcam_height,
                           int gl_size, unsigned char thresh)
{
    unsigned char max_delta=0;
    int i=0, j=0;
    int length;
    unsigned char *temp_delta = (unsigned char *)malloc(
        sizeof(unsigned char)* webcam_width * webcam_height);

    length = webcam_width * webcam_height;

    webcam_tex.which_buffer = webcam_tex.which_buffer ? 0 : 1;

    if(!video_hdr->lpData)
        return FS_TRUE;

    for(i=0; i<length; i++)
    {
        // Save the current frames data for comparison on the next frame
        // NOTE: Were only comparing the red channel (lpData is BGR), so
        // in theory if the user was in a solid red room, coated in red
        // paint, we wouldn't detect any movement....chances are this
        // isn't the case :) For our purposes, it this test works fine
        webcam_tex.back_buffer[webcam_tex.which_buffer][i]
        = video_hdr->lpData[i*3];

        // Compute the delta buffer from the last frame
        // If it's the first frame, it shouldn't blow up given that we
        // cleared it to zero upon initialization
        temp_delta[i] =
        abs(webcam_tex.back_buffer[webcam_tex.which_buffer][i] -
```

```
    webcam_tex.back_buffer[!webcam_tex.which_buffer][i]);

    // Is the difference here greater than our threshold?
    if(temp_delta[i] > max_delta)
        max_delta = temp_delta[i];
}

// Fit to be inside a power of 2 texture
for(i=0; i<webcam_height; i++)
{
    memcpy(&delta_buff[i*(gl_size)],
           &temp_delta[i*(webcam_width)],
           sizeof(unsigned char)*webcam_width);
}

free(temp_delta);

if(max_delta > thresh)
    return TRUE;
else
    return FALSE;
}
```

### 1.22.2 操纵 Web 摄像机数据

#### 1. 取得 BGR 像素

完成所有的测试和收集工作后，便可以对 Windows 发送来的数据进行操纵。我们只是将 VIDEOHDR 结构中的像素（Windows 返回的本机格式为 BGR）复制到前面分配的、大小为 2 的整数次幂的缓冲区中。注意，这种技术避免了调整纹理数据的像素大小，因为它直接复制像素，保留高宽比不变。这种技术的惟一缺点是，将导致纹理中有空白，图像的最上面将是一条黑带。可以通过操纵纹理坐标（作为贴图附加到 3D 几何体上后）或调整纹理的大小来消除该黑带。

```
// WEBCAM_MAKE_BGR
void webcam_make_bgr(unsigned char *bgr_tex, unsigned char *vid_data, int webcam_width,
int webcam_height, int glsize)
{
    int i;

    for(i=0; i<webcam_height; i++)
    {
        memcpy(&bgr_tex[i*(glsize*3)],
               &vid_data[i*(webcam_width*3)],
               sizeof(unsigned char)*webcam_width*3);
    }
}
```

## 2. 转换为灰度值

捕获 BGR 数据后, 可将其转换为灰度值。这种转换将导致图像的大小为纹理的 1/3。对于 Internet 连接的速度很慢, 但仍想传输 Web 摄像机数据的用户来说, 这很有用。

下面的函数将颜色缓冲区中的 RGB 组分乘以一个缩放系数, 这实际上将三个颜色通道压缩为一个。

```
// WEBCAM_MAKE_GREYSCALE
void webcam_make_greyscale(unsigned char *grey,
                           unsigned char *color, int dim)
{
    int i, j;

    // Greyscale = RED * 0.3f + GREEN * 0.4f + BLUE * 0.3f
    for(i=0, j=0; j<dim*dim; i+=3, j++)
    {
        grey[j] = (unsigned char)float_to_int(0.30f * color[i] +
                                              0.40f * color[i+1] +
                                              0.30f * color[i+2]);
    }
}
```

## 3. 展示现实生活的卡通

将所有数据转换为灰度值后, 便可以通过操纵这些数据, 以类似卡通的方式来绘制图片。这种方法将图像分为 5 种不同的等级和 6 种不同的颜色, 并使用固定的值给不同范围的像素值着色。我们所需做的只是做一些简单的比较, 并根据热强度常量对每个像素做出判断。

我们将最终的结果同来自灰度缓冲区或增量缓冲区的对照表进行比较。如果每帧都要看到该图像(单缓冲区), 则需要同灰度进行比较。为获得不同的结果, 我们将根据热强度常量给每个像素指定一个随机的色强。

```
// WEBCAM_INIT_CARTOON
void webcam_init_cartoon(cartoon_s *cartoon_tex)
{
    char i;

    for(i=0; i<3; i++)
    {
        // Pick random colors in our range
        cartoon_tex->bot_toll_col[i] = rand()%255;
        cartoon_tex->min_toll_col[i] = rand()%255;
        cartoon_tex->low_toll_col[i] = rand()%255;
        cartoon_tex->med_toll_col[i] = rand()%255;
        cartoon_tex->high_toll_col[i] = rand()%255;
        cartoon_tex->max_toll_col[i] = rand()%255;
    }
}
```

```
#define MIN_CAM_HEAT 50
#define LOW_CAM_HEAT 75
#define MED_CAM_HEAT 100
#define HIGH_CAM_HEAT 125
#define MAX_CAM_HEAT 150

// WEBCAM_MAKE_CARTOON
void webcam_make_cartoon(unsigned char *cartoon,
                        cartoon_s cartoon_tex,
                        unsigned char *data, int dim)
{
    int i, j, n;

    for(i=0, j=0; j<dim*dim; i+=3, j++)
    {
        if(data[j] < MIN_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.bot_toll_col[n];
        }
        if(data[j] > MIN_CAM_HEAT && data[j] < LOW_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.min_toll_col[n];
        }
        if(data[j] > LOW_CAM_HEAT && data[j] < MED_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.low_toll_col[n];
        }
        if(data[j] > MED_CAM_HEAT && data[j] < HIGH_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.med_toll_col[n];
        }
        if(data[j] > HIGH_CAM_HEAT && data[j] < MAX_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.high_toll_col[n];
        }
        if(data[j] > MAX_CAM_HEAT)
        {
            for(n=0; n<3; n++)
                cartoon[i+n] = cartoon_tex.max_toll_col[n];
        }
    }
}
```

#### 4. 上载新的纹理

现在，余下的工作只有将纹理上载到 OpenGL。首先，要从 Video for Windows 那里获得色值。计算出新的色值后，便可以将其转换为灰度值，然后实现卡通渲染。完成所有的图像操纵工作后，我们调用 `glTexSubImage2D()` 来将其转换为合适的纹理，这样 3D 应用程序便可以使用它了。

```
// WEBCAM_MAKE_TEXTURE
void webcam_make_texture(LPVIDEOHDR video, webcam_draw_mode mode)
{
    // Build the color first
    webcam_make_bgr(webcam_tex.bgr,
video->lpData,
wco_cam_width,
wco_cam_height,
webcam_tex.size);

    if(mode == GREYSCALE || mode == CARTOON)
        webcam_make_greyscale(webcam_tex.greyscale,
                                webcam_tex.bgr, webcam_tex.size);

    // Note: Could also pass in the delta buffer instead of
    // the greyscale
    if(mode == CARTOON)
        webcam_make_cartoon(webcam_tex.bgr,
                                webcam_tex.cartoon,
                                webcam_tex.greyscale,
                                webcam_tex.size);

    //////////////////////////////////////

    // Upload the greyscale version to OpenGL
    if(mode == GREYSCALE)
    {
        glBindTexture(GL_TEXTURE_2D, webcam_tex.gl_grey);
        glTexSubImage2D(GL_TEXTURE_2D, 0, 0, 0,
                                webcam_tex.size, webcam_tex.size,
                                GL_LUMINANCE,
                                GL_UNSIGNED_BYTE, webcam_tex.greyscale);
    }
    // Upload the color version to OpenGL
    else
    {
        glBindTexture(GL_TEXTURE_2D, webcam_tex.gl_bgr);
        glTexSubImage2D(GL_TEXTURE_2D, 0, 0, 0,
                                webcam_tex.size, webcam_tex.size,
                                GL_BGR_EXT, GL_UNSIGNED_BYTE, webcam_tex.bgr);
    }
}
```



## 5. 释放 Web 摄像机窗口

使用完 Web 摄像机后，需要释放该窗口，并将回调函数设置为 NULL，这样 Windows 便会停止向它发送消息。另外，还必须释放前面给颜色缓冲区、灰度缓冲区和增量缓冲区分配的内存。

```
// WEBCAM_DESTROY
void webcam_destroy(void)
{
    if(cam_driver_on)
    {
        capSetCallbackOnFrame(hWndCam, NULL);
        DestroyWindow(hWndCam);
        hWndCam = NULL;

        if(webcam_tex.bgr)
            free(webcam_tex.bgr);

        if(webcam_tex.greyscale)
            free(webcam_tex.greyscale);

        if(webcam_tex.delta_buffer)
            free(webcam_tex.delta_buffer);

        if(webcam_tex.back_buffer[0])
            free(webcam_tex.back_buffer[0]);

        if(webcam_tex.back_buffer[1])
            free(webcam_tex.back_buffer[1]);
    }
}
```

### 1.22.3 结论

Web 摄像机还有很多游戏开发人员没有挖掘出来的潜力。通过跟踪颜色对象，并将其旋转数据从 2D 转换为 3D[Wu99]，它可用作输入设备（就像使用鼠标一样）。使用 Web 摄像机，通过对输入帧执行数据修匀和颜色跟踪算法，它甚至可以替代标准鼠标。

### 1.22.4 参考文献

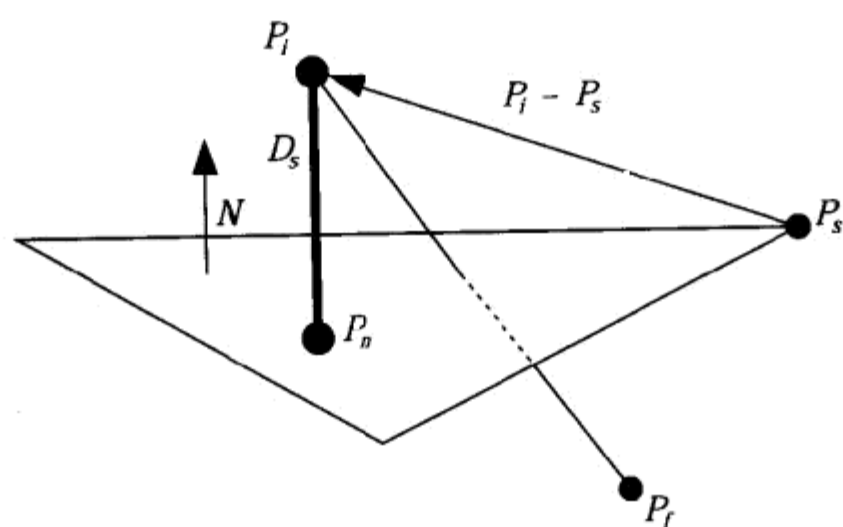
Microsoft Developer Network Library <http://msdn.microsoft.com/library/devprods/vs6/visualc/vcsample/vcsmcptest.htm>.

[Firetoad00] Firetoad Software, Inc., *Ectosaver*, 2000 [www.firetoads.com](http://www.firetoads.com).

[Wu99] Wu, Andrew, "Computer Vision REU 99" [www.cs.ucf.edu/~vision/reu99/profile-awu.html](http://www.cs.ucf.edu/~vision/reu99/profile-awu.html).



# 数学技巧





## 绪 论

Eddie Edwards, Naughty Dog

eddie@tinyted.net

长期以来，计算机和数学一直相伴相生。作者最早的记忆之一是见到 Cray 超级计算机。在 20 世纪 80 年代的英国，这种计算机试图被用来预测该国不可预知的天气情况。这种计算机体积庞大——占据一个小房间，每秒执行数百万次浮点运算。那时，我那基于 6809 的微机却在为实时显示线框立方体而奋斗！当时，一般的程序员并不会将很多时间花在考虑数学方面的问题上，而操纵这种庞然大物的程序员正忙于采用任何游戏程序员都将引以为荣的技巧解决一些非常棘手的问题。

到 2001 年，人们起居室中的游戏控制台的计算能力已远远超过 Cray XMP，而这些小型机器的程序员正是本书的读者。然而，我们仍然无法预知伦敦明天是否会下雨！

本章阐述了这种进展过程。一方面，Yossarian King 撰写的关于浮点计算技巧的文章介绍了通过现代浮点技术赋予传统的游戏编程技巧以现代意义；另一方面，Alex Vlachos 和 John Isidoro 撰写的文章使用一些复杂的数学知识，推导出一种平滑的、低成本的四元数插值方案。Steve Rabin 介绍的空间快速分组方法提醒我们，不管我们拥有的计算能力有多强大，有些问题仍然很困难（虽然由于他提供的算法，问题解决起来会容易些）。

这些文章同时展示了游戏编程的另一面。这里的每篇文章都展示了作者高度的创造力。例如，数学知识丰富的读者可能无视 Jesse Laeuchli 对不规则碎片的介绍；然而，这样他们将错过 Jesse 的重大发明——一个旨在提高游戏速度的噪声函数。

其他的 4 篇文章对读者同样很有用。Aaron Nicholls 解决了一系列有关导弹发射的问题，这对任何 3D 游戏都很有帮助；John Olsen 对矢量和平面方面的技巧进行了总结，提供了最常见的矢量计算方程；Graham Rhodes 提供了一种线段相交问题的解决方案；而 Carl Dougan 描述了一种简单但高度稳定的曲线跟踪（3D 问题沿曲线前行）算法。

鉴于上述因素，编写本章令人非常高兴。希望读者阅读后同我一样感到满意。



## 2.1 浮点计算技巧：使用 IEEE 浮点格式以提高性能

Yossarian King, Electronic Arts Canada  
yking@ea.com

### 2.1.1 概述

**整**数的精度和量级都是固定的，而浮点数的小数点位置和量级都不固定。历史上，整数的计算速度高于浮点数，因此游戏程序员使用整数，而避免使用浮点数。对于一般的计算，使用整数时比较繁琐，但由于其性能方面的优势，这样做是值得的。然而，硬件的价格已经降低，当前的 PC 只需几个周期就能够完成浮点数的加减乘除运算。现在，使用浮点数运算很容易，游戏程序员可以充分利用这一点。

虽然基本的浮点算术的速度很快，但复杂的函数仍然很慢。可以对浮点函数库进行优化，但实现这些函数通常旨在提高精度而不是性能。就游戏而言，性能常常比精度更重要。

本节介绍各种提高浮点运算性能的技巧，这些技巧以降低精度为代价来提高执行速度。很长时间以来，查表（table lookup）一直是一种用于整数数学计算的标准技巧，本节将介绍通用的线性和对数查表技巧，用于优化任何浮点函数。

接下来的几个小节将讨论以下内容：

- IEEE 浮点标准；
- 对浮点数/整数进行快速转换、比较和定位（clamping）的技巧；
- 用于优化正弦和余弦函数的线性查表法；
- 用于优化平方根函数的对数查表法；
- 用于优化任何浮点函数的通用查表法；
- 性能测量的重要性。

### 2.1.2 IEEE 浮点格式

IEEE 浮点数标准规定了一种二进制表示，同时定义了关于四舍五入、精度和异常结果（如被零除）的规则。本节介绍的技巧依赖于二进制表示，但通常不关心四舍五入和异常处理。如果计算机或游戏控制台使用这种标准的二进制表示，则无论浮点处理是否完全遵循 IEEE 标准，这些技巧都适用。Pentium III SSE 和 PS2 矢量单元（vector unit）实现了 IEEE 标准中不支持异常处理的子集，但这里介绍的技巧适用于这些指令集。

在 IEEE 标准中，使用一个符号位、一个偏置指数和一个规格化尾数来表示浮点数。32 位的单精度浮点数（C 语言中的 float）按如图 2.1.1 所示的格式存储。



图 2.1.1 在 IEEE 标准中，32 位的浮点数用符号（1 位）、指数（8 位）和尾数（23 位）表示

指数按偏置格式被存储为一个正数（而不是整数中使用的 2 的反码表示），存储的值为实际指数与 127 的和；尾数通常以规格化格式存储，这意味着 23 位表示的是小数部分，整数部分为 1。以这样的方式规格化可使精度最高。

因此，浮点数由一个表示 1~2 的规格化尾数、一个表示二进制小数点位置的偏置指数和一个符号位组成。表示的浮点数值为：

$$n = s \times 1.m \times 2^{(e-127)}$$

例如，-6.25 的二进制表示为 -110.01 或  $-1 \times 1.1001 \times 2^2$ ，因此， $s = 1$ ， $e = 2 + 127 = 10000001$ ， $m = [1.]1001$ ，如图 2.1.2 所示。

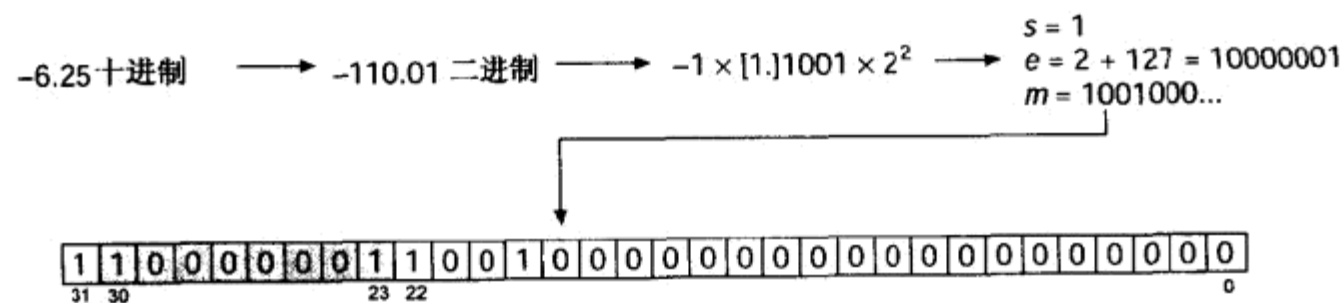


图 2.1.2 使用 32 位的 IEEE 浮点格式时，-6.25 被存储在内存中的情况

另外，我们将使用特定的指数值来表示一些“不可思议的值”。当  $e$  为 255 时， $m$  指示诸如不是数字（NaN）、未定义的结果或正（负）无穷大等特殊情况。对于非规格化数字——非常小，以至于指数无法用 8 位表示出来，则  $e$  为 0。

64 位的双精度浮点数的存储格式与此相同，只是指数占 11 位，而尾数占 52 位。计算实际的指数时，将减去 1023，而不是 127。双精度数占用的空间为单精度数的两倍，从内存中装载它们以便进行处理时速度可能更慢。因此，在游戏编码中，通常应避免使用双精度数。本节只使用单精度浮点数。

2.1.3 浮点数技巧

介绍查表技巧之前，先讨论一些浮点数技巧，以帮助解释对浮点數位模式的处理技巧。



## 1. 浮点数/整数转换

接下来介绍的查表技巧将浮点数转换为整数以创建一个查表索引。在 Pentium II 中，这种操作的速度可能很慢。例如，使用 `(int)f` 将浮点数强制转换为整数需要 60 个时钟周期。这是因为 ANSI C 标准规定，将浮点数转换为整数时，应截除小数部分，而在默认情况下，FPU 将浮点数四舍五入到最接近的整数。因此，将浮点数转换为整数将成为一个函数调用，被调用的例程修改 FPU 舍入模式、执行转换然后恢复舍入模式。

在整数和浮点数之间进行转换的成本取决于您使用的编译器和处理器。和所有的优化一样，应将这种转换技巧同常规类型转换进行比较，并将代码进行反汇编，看看实际的情况如何。

可以将浮点数加上  $1 \times 2^{23}$ ，然后丢弃结果中高指数位，这样转换速度将快得多。我们首先来看代码，然后分析其中的原理。

为此，定义下面的共用体将很有帮助，它让我们能够存取 32 位的浮点数或整数。

```
typedef union
{
    int      i;
    float    f;
} _INTORFLOAT;
```

本节的代码片段都将使用类型 `INTORFLOAT`。它使得存取数字的位模式非常简单——实际上，编译器生成的代码可能比您想象的要多得多。例如，在 Pentium II 上，浮点数寄存器和整数寄存器位于不同的硬件中，数据可能从一个寄存器被移到另一个寄存器，而无需经过内存，因此存取共用体 `INTORFLOAT` 的成员时，可能需要装载其他内存和存储器。

将浮点数转换为整数的代码如下：

```
INTORFLOAT  n;          // floating-point number to convert
INTORFLOAT  bias;       // "magic" number

bias.i = (23 + 127) << 23; // bias constant =  $1 \times 2^{23}$ 
n.f = 123.456f;           // some floating-point number

n.f += bias.f;            // add as floating-point
n.i -= bias.i;            // subtract as integer
// n.i is now 123 - the integer portion of the original n.f
```

为何这样做管用呢？加上浮点数  $1 \times 2^{23}$  后，尾数将被存储在右边的 23 位中，同时指数将被设置为一个已知的值—— $23 + 127$ ；通过整数方式减去已知的指数后，清除了结果中高位值，只留下低位值。图 2.1.3 以 43.25 为例说明了这一过程。

在 Pentium II（所有的东西都在高速缓存中）上，这将转换时间从 60 个时钟周期缩短为大约 5 个。也可以通过编写内联汇编代码，让 FPU 在将浮点数转换为整数时不改变舍入模式，这样速度将比强制类型转换快，但通常比这里介绍的技巧慢。

只要被转换的浮点数不与偏置常数“重叠”，这种技巧都管用。换句话说，只要浮点数小于  $2^{23}$ ，这种技巧都管用。





```

...
INTORFLOAT atmp, btmp;
atmp.f = f; btmp.f = b;
if (atmp.i < btmp.i) // integer compare
...

```

通常，整数比较的流水线技术更佳，速度更快。注意，一种例外情况是  $a$  和  $b$  都为负数，因为指数和尾数并非像整数比较期望的那样以 2 的补码方式存储的。如果至少有一个数为正，则采用上述方式进行比较的速度将更快。

#### 4. 定位 (clamping)

在游戏编程中，经常需要将一个值限定在特定的范围内，通常为  $[0, 1]$ 。通过将符号位转换为掩码 (mask)，可以将浮点数  $f$  转换为 0 (即如果  $f < 0$ ，则将  $f$  设置为 0)，如下面的代码片段所示：

```

INTORFLOAT ftmp;
ftmp.f = f;
int s = ftmp.i >> 31; // create sign bit mask
s = ~s;                // flip bits in mask
ftmp.i &= s;            // ftmp = ftmp & mask
f = ftmp.f;

```

对  $f$  执行右移 31 位的移位运算——将所有 32 位的值都设置为符号位的值，然后将结果赋给  $s$ 。对  $s$  执行非运算将创建一个所有位都为 0 (如果  $f$  为负) 或 1 (如果  $f$  为正) 的值。对  $s$  和  $f$  执行与运算将使  $f$  的值不变或变为 0。最终的结果是：如果  $f$  为负，则变为 0；如果  $f$  为正，则不变。

上述代码完全是在整数单元中运行的，而且既没有比较，也没有分支。在测试代码中，浮点数比较和定位需要大约 18 个时钟周期，而整数定位少于 5 个时钟周期 (这些时间中包含循环的开销)。

将正数定位为 0 (如果  $f > 0$ ，则将其设置为 0) 的用处不大，但更容易，因为此时不需要反转掩码中的各个位。

```

INTORFLOAT ftmp;
ftmp.f = f;
int s = ftmp.i >> 31; // create sign bit mask
ftmp.i &= s;            // ftmp = ftmp & mask
f = ftmp.f;

```

要定位到 1，可减去 1、定位到 0，然后再加上 1。

```

INTORFLOAT ftmp;
ftmp.f = f - 1.0f;
int s = ftmp.i >> 31; // create sign bit mask
ftmp.i &= s;            // ftmp = ftmp & mask
f = ftmp.f + 1.0f;

```

注意，在汇编语言中，使用条件装载指令通常可提高定位运算的速度，因为这样可消除分支，从而消除指令管道中的分支预测逻辑。

### 5. 绝对值

这很容易，因为浮点数不使用 2 的补码，所以要求浮点数的绝对值，只需将符号位设置为 0 即可。

```
INTORFLOAT ftmp;
ftmp.f = f;
ftmp.i &= 0x7fffffff;
f = ftmp.f;
```

注意，与对 f 求反之前通过比较来判断 f 是否小于 0 相比，上述代码的速度要快得多。

## 2.1.4 用于正弦和余弦函数的线性查找表

在游戏中，三角学常常很有用——用于计算距离和角度、沿圆圈前进、制作水网（water mesh）动画。标准数学库中包含各种三角函数，但它们的速度慢，使用的是双精度值，因此占用的内存比需要的多。在游戏中，低精度计算通常就足够了。

为高效地计算正弦和余弦，可使用一个查找表。一种常见的方法是使用定点数学，其中使用某个范围（如 0~1023）内的整数来表示角度，这个范围覆盖了整个圆周。然而，这意味着游戏程序员必须理解正弦和余弦的库函数实现，并按其要求的格式来表示角度。通过使用浮点技巧来实现高效的索引，可以创建使用弧度的浮点三角函数，且不要求程序员了解实现的细节。

### 1. 正弦函数

我们来实现下述函数：

```
float fsin( float theta );
```

这可以通过一个查找表来轻松地实现。采用下述方式初始化一个包含 256 个条目的表，它覆盖了  $0\sim 2\pi$  的角度：

```
sintable[i] = (float)sin((double)i * 2.0*3.14159265/256.0)
```

上述代码将 i（取值范围为 0~256）转换为一个浮点数弧度（ $0\sim 2\pi$ ），然后计算该角度的正弦。

有了上述表后，可以这样实现函数 fsin：

```
float fsin( float theta )
{
    i = (unsigned int)(theta * 256.0f/(2.0f*3.14159265f));
    return table[i];
}
```

然而，这存在两个问题。首先，它采用了速度较慢的浮点数到整数的类型转换；其次，

如果  $\theta$  不在范围  $[0, 2\pi]$  之内，则索引将不在表内。

下述实现解决了这两个问题：

```
#define FTOIBIAS 12582912.0f // 1.5 * 2^23
#define PI      3.14159265f

float fsin( float theta )
{
    int i;
    INTORFLOAT ftmp;
    ftmp.f = theta * (256.0f/(2.0f*PI)) + FTOIBIAS;
    i = ftmp.i & 255;
    return table[i];
}
```

该实现使用前面介绍的浮点数偏置技巧来完成浮点数到整数的快速转换。它将整数同 255 进行与运算，从而确保索引值总是位于 0~255 的范围之内。注意，如果  $f$  的值大于  $2^{22}$ ，上述浮点数到整数的转换技巧将无效，因此，必须定期地将  $f$  的值减小到有效范围  $[0, 2\pi]$  之内。

在 Pentium II 上，上述  $\text{fsin}$  实现需要 10 个时钟周期（假设所有数据和代码都位于主高速缓存中），而标准数学库中的  $\sin$  实现需要 140 个时钟周期（虽然  $\sin$  使用了 FPU 中的硬件正弦指令）。



包含 256 个条目的浮点数表占用大约 1K 的空间，在内循环期间，很容易将其保留在高速缓存中。准确度为 8 位，这是由表的大小决定的。通过分析查找表，可以轻松地检测到最坏情况下的错误（附带光盘中的代码演示了这一点）。增大查找表可以提高精度，但将降低高速缓存的性能。

## 2. 余弦函数

余弦函数可以采用同样的方式来实现，并使用单独的查找表，但我们可以利用这样一个事实，即  $\cos(\theta) = \sin(\theta + \pi/2)$ ，从而使用同一个查找表。为此，只需将索引加上 256/4（因为加上  $\pi/2$  相当于加上圆周的 1/4），这可以在偏置指数时完成。余弦函数的实现代码如下：

```
float fcos( float theta )
{
    int i;
    INTORFLOAT ftmp;
    ftmp.f = theta * (256.0f/(2.0f*PI)) + FTOIBIAS + 64f;
    i = ftmp.i & 255;
    return table[i];
}
```

根据应用，同时实现正弦函数和余弦函数通常很有用。与分别计算相比，这样做的效率更高——查找正弦值，然后将索引值加上 64 并将结果与 255 执行与运算，以查找余弦值。如果需要同时计算多个正弦和余弦值，可编写自定义的代码，交替进行计算，从而提高速度。

### 2.1.5 平方根函数的对数优化

在游戏中，平方根对计算距离、矢量规格化和解二次方程等很有用。虽然 FPU 中内置了平方根指令，但在 Pentium II CPU 中，标准 C 函数库中的 `sqrt` 函数仍需要大约 80 个时钟周期，因此它是另一个适合优化的函数。

平方根优化是通过操纵浮点数位来实现的。由于平方根的对数和通用换算特性，我们可以将平方根计算进行分解，从而分别操纵尾数和指数。来看浮点数的平方根：

$$\begin{aligned} f &= 1.m \times 2^e \\ \text{sqrt}(f) &= \text{sqrt}(1.m \times 2^e) \\ &= \text{sqrt}(1.m) \times 2^{e/2} \end{aligned}$$

因此，要计算  $f$  的平方根，可分别计算尾数的平方根和指数的一半。然而，指数为整数，因此如果指数为奇数，则将其除以 2 后，最后一位将丢失。为解决这种问题，可将最低一位放到尾数中，即：

$$\text{sqrt}(f) = \text{sqrt}(1.m \times 2^{e_0}) \times 2^{\lfloor e/2 \rfloor}$$

其中  $e_0$  为指数的最后一位。

因此，可以使用一个包含 256 个条目的尾数平方根表，并调整指数的计算方式来实现平方根函数，如下所示：

```
float fsqrt( float f )
{
    INTORFLOAT      ftmp;
    unsigned int     n, e;

    ftmp.f = f;
    n = ftmp.i;
    e = (n >> 1) & 0x3f800000;    // divide exponent by 2
    n = (n >> 16) & 0xff;        // table index is e0+m22-m16

    ftmp.i = sqrttable[n] + e;    // combine results

    return ftmp.f;
}
```

表索引为尾数的左边 7 位和指数的最后一位。查找表中包含尾数的平方根。

为计算平方根的指数，将  $f$  的指数向右移一位——将其除以 2。由于指数被偏置，这会把偏置常量和指数都除以 2，而这并非我们希望的。为此，将 `sqrttable` 条目与另一个因子相加，以重新偏置指数。

在 Pentium II CPU 上，上述 `fsqrt` 函数需要 16 个时钟周期——比 C 函数库实现快 5 倍。同样，这里也假设数据和代码都位于高速缓存中。





附带光盘中的代码更详细地解释了这种算法。

### 2.1.6 优化任何函数

考虑一个任意的接受一个变量的函数：

$$y = f(x)$$

前面介绍的技巧揭示了两种基于表来优化一般性函数的基本方法。对于正弦和余弦函数， $x$  的值在已知的范围内被线性地量化，并被用作查找  $y$  值的表索引。对于平方根函数， $x$  的值以对数的方式被量化，并被用作查找一个值的表索引，然后将找到的值乘以一个以  $x$  的指数为自变量的函数，得到  $y$  的值。

线性方法通过线性量化，将浮点数按一定的比例缩放，并将其转换为一个整数，通过线性量化生成查找表索引。这种技巧很简单，与整数查找表极其类似，其中惟一的技巧是高效地将浮点数转换为整数索引。对数方法直接使用浮点数的位模式作为索引，以实现对数量化。

这两种技巧都可推广用于优化任何函数。到底是线性方法合适还是对数方法合适，取决于要优化的函数。

#### 1. 线性量化

前面的 `fsin` 函数可用作通过线性量化优化函数的模板。假设函数只用于这样的范围： $x \in [A, B]$ ，则可以构建一个均匀地覆盖该范围的查找表，并计算该范围内某个特定的  $x$  值对应的表索引。因此，优化后的  $f$  函数的实现如下：

```
#define FTOIBIAS      12582912.0f  // 1.5 * 2^23
#define TABLESIZE    256
#define INDEXSCALE    ((float)TABLESIZE / (B - A))

float  flut( float x )
{
    int      i;
    INTORFLOAT  ftmp;
    ftmp.f = x * INDEXSCALE + (FTOIBIAS - A * INDEXSCALE);
    i = ftmp.i & (TABLESIZE - 1);
    return ftable[i];
}
```

初始化查找表的代码如下：

```
ftable[i] = f( (float)i / INDEXSCALE + A );
```

其中  $f$  是该函数精确的浮点实现。`flut` 需要执行两种浮点运算（加和乘）、一个整数按位掩码和一次查表。在 Pentium II CPU 上，这需要大约 10 个时钟周期。



可以通过在两个最接近的表项之间进行线性插值来提高精确度, 但这样需要的时钟周期将多几个。附带光盘提供了一个支持这种函数优化的 API, 其中包含可选的、用于提高精确度的线性插值。

## 2. 对数量化

线性方法在范围[A, B]之间均匀地进行量化。根据要优化的函数, 采用对数方法可能更合适, 如平方根函数的优化。这里的基本思想是, 直接将浮点数中的位用作查找表索引, 而不是将浮点数转换到某个整数范围内。通过提取选中的符号位、指数位和尾数位, 可以将 1:8:23 的 IEEE 浮点数转换为简化的格式——根据我们的爱好, 使用不同的位数来表示符号、指数和尾数。

在平方根范例中, 我们提取了 8 位, 实现一种 0:1:7 的对数量化表示。我们使用了指数中的 1 位和尾数中的 7 位, 而忽略了符号位, 因为负数的平方根没有定义。0:1:7 格式表示一个 8 位的尾数(别忘了 IEEE 表示法中的 1)和一个 1 位的指数, 因此其取值范围为  $[1].0000000 \times 2^0 \sim [1].1111111 \times 2^1$ , 即 1~4。

平方根函数被分解为对 0:1:7 量化数字的运算(表索引)和对指数的运算(除以 2)。在优化这两种运算并将尾数和指数合并为一个 32 位的浮点数时, 还采用了其他技巧。

这种对数量化方法也可用于优化其他函数。IEEE 格式使得只需执行移位和与运算, 便可以轻松地提取指数的最后几位和尾数的最左边几位。要提取指数的 *ebits* 位和尾数的 *mbits* 位, 只需这样做:

$$bits = (n \gg (23 - mbits)) \& ((1 \ll (ebits + mbits)) - 1)$$

上述代码将数字 *n* 向右移位, 使得所需的指数位和尾数位位于最右边, 然后屏蔽所需的位数。

根据要优化的函数, 可通过位运算来处理符号位。如果确定处理的只能是整数(如计算平方根)或函数总是返回一个正值, 则可以忽略符号位。如果结果的符号总是与输入值的符号相同(即  $f(-x) = -f(x)$ ), 则只需存储并恢复符号位即可。

对于输入值在特定范围内的函数, 屏蔽选中的指数位和尾数位可获得一个表索引。例如, 如果您只关心范围[1,16]的情况, 则可以使用 2 个指数位和 4 个尾数位。这种 0:2:4 的表示法能够存储的二进制值为  $1.0000 \times 2^0 \sim 1.1111 \times 2^3$ , 相应的十进制值为 1.0~15.5。将这些位屏蔽掉, 并将它们直接用作一个预先计算好的包含 64 个条目的查找表的索引, 这需要的时钟周期很少, 计算速度很快。然而, 如果要提高精度, 表将增大, 而且可能非常大, 这种情况下, 高速缓存的性能将降低。

另一种方法是将指数计算和尾数计算分开, 就像平方根范例中那样。如果函数  $f(x)$  可被分解为:

$$f(x) = f(1.m \times 2^e) = f1(1.m) \times 2^{f2(e)}$$

则可以使用尾数 *m* 中的 8 位创建一个包含 256 个条目的查找表, 以近似计算 *f1* 的值, 并直接计算 *f2*——对指数 *e* 的整型运算。这实际上是前面的平方根范例中使用的技巧。

对数量化是一种功能强大的工具, 但要求的位运算随要优化的函数而异。技巧并非总是

完全通用的，但这里介绍的方法对处理特定的优化问题肯定会有所帮助。

### 2.1.7 性能测量

优化代码时，务必要在优化之前和之后仔细地测量性能。由于高速缓存的行为、对分支的错误预测或编译器的糟糕处理，有时候理论上很不错的优化在实现时可能带来麻烦。进行修改时，一定要确保性能得到了提高——决不要想当然。

确保启用了编译器的优化功能。在合适的情况下，使用内联函数 (Inline function)。同样，使用内联函数或调整编译器设置后，应对结果进行测试。

对代码进行比较时，确保编译器优化没有妨碍测试。对代码进行反汇编，并单步执行它，以确保它按您期望的方式运行。进行计时时，循环运行一些代码常常很有帮助——但如果循环被优化掉，则计时将不会非常准确。

在 Pentium 计算机上，可以使用 rdtsc (读取时间戳计数器) 指令来获得当前的 CPU 周期数。Intel 警告说，使用结果之前，必须执行该指令多次。Intel 还建议使用诸如 cpuid 等将刷新指令高速缓冲的指令，这样将获得更一致的计时结果。要获得绝对时间，可将周期数除以处理器的执行速度 (MHz)，将其转换为秒数。

周期计数器是最可靠的测量细粒度性能的方式。对于高级性能剖析，诸如 VTune 和 TrueTime (用于 PC) 等工具很有用。对于任何基准点确定，都必须确保内存行为是真实的，因为内存瓶颈是影响当代处理器性能的最重要的因素之一。一定要知道基准是使用高速缓存的，并尽可能地模拟游戏的高速缓存行为。对于基准，可在计时之前运行算法多次，给高速缓存“热身”。然而，热身后的高速缓存可能未能模仿游戏的行为——最好的办法是直接在游戏中建立基准。禁用中断可获得更可靠的结果，或测量多次并忽略峰值。



本节中指出的所有周期数都是基于 Intel Pentium II 450-MHz 机器的。每种运算都在一个循环中重复 1000 次，并测试了多次，给指令和数据高速缓存热身。周期数包含循环开销。有关实际使用的基准，请参阅附带光盘中的代码。

如果查找表一直保留在高速缓存中，则本节中介绍的查找表技术是合适的。在渲染管道和物理引擎的内循环中，情况可能确实如此；但如果您在代码中随机地调用这些函数，则情况可能并非如此。如果查找表没有保留在高速缓存中，则进行更多计算并更少存取内存的技巧可能更合适——如多项式逼近 (参见[Edwards00])。

### 2.1.8 结论

本文只是对浮点优化做了蜻蜓点水式的介绍。这里探讨的主要是查找表方法，这种方法常常能极大地提高速度。然而，一定要了解高速缓存的行为，并对结果进行基准化。有时候，采用更多地进行计算但更少地存取内存的方法——如多项式逼近，能够获得更高的速度。这里介绍的技巧已被通过各种方式进行扩展，同时还有许多其他的技巧。正如一本流行的图书指出的，代码优化艺术是一种禅宗，像本文这样简要的介绍是无法涵盖其全部内容的。

### 2.1.9 参考文献

---

[Abrash94] Abrash, Michael, *Zen of Code Optimization*, Coriolis Group, 1994.

[Edwards00] Edwards, Eddie, "Polynomial Approximations to Trigonometric Functions," *Game Programming Gems*, Charles River Media, 2000.

[Intel01] Intel Web page on floating-point unit and FPU data format. Good for PCs, but relevant to any IEEE-compliant architecture. Available at <http://developer.intel.com/design/intarch/techinfo/Pentium/fpu.htm>.

[Lalonde98] Lalonde, Paul, and Dawson, Robert, "A High Speed, Low Precision Square Root," *Graphics Gems*, Academic Press, 1998.



## 2.2 矢量和平面技巧

John Olsen, Microsoft

infix@xmission.com

**您**的碰撞检测 (collision detection) 函数正常运行, 当您向它传递一个位置和一个速度向量时, 它返回曲面上的一个点和一条法线。实际上, 根据已有的数据, 您能够做的工作很多。

通常, 您可能希望碰撞函数返回实际的碰撞点, 但本文的方法演示了如何处理这样的碰撞结果, 即它只指出了交会平面。由于平面可以通过平面上的任何一个点和法线来指定, 因此可以通过数学处理来获得所有需要的信息。

提供给碰撞函数的数据起点  $P_i$  和终点  $P_f$ , 如果发生碰撞, 则输出为一个平面, 该平面是由单位向量面法线  $N$  和面上的一个点  $P_s$  定义的。该点不必是实际的交点, 而只需是平面上的一个点即可。

出于优化的目的, 您可能想将这些计算结果的一部分放到碰撞代码中。检测时, 所需的大部分信息都已计算好。与重新计算相比, 重用已知的信息将快得多。

平面方程  $Ax + By + Cz + D = 0$  映射提供的数据。其中  $x$ 、 $y$ 、 $z$  是向量  $N$  的分量, 而  $D$  是点积  $N \cdot P_s$ 。

### 2.2.1 相对于碰撞面的高度

检测碰撞时, 最常用的数据之一是数据点的高度。如果高度为正, 则该点位于平面的上方, 碰撞还未发生; 如果为负, 则已经发生碰撞, 并已穿过平面。

如果检测点位于检测面的两边, 则典型的碰撞检测代码将指出发生了碰撞。这意味着如果要对碰撞进行预测, 需要传递一个位置和一个被放大的速度向量。这样与实际的移动情况相比, 被放大的向量将更早地与平面相交。

让碰撞代码返回平面上的一个点和法线后, 便可以使用起点位置获得相对于该平面的高度。计算高度时, 没有用到终点位置。

如图 2.2.1 所示, 我们要计算投影到面法线  $N$  上时, 向量  $(P_s - P_i)$  的长度, 这样便可以知道从点到面的最短距离。根据定义, 该最短向量是到平面的垂线。这可以通过点积计算得到, 因此到面  $D_s$  (非向量) 距离为:

$$D_s = (P_i - P_s) \cdot N$$

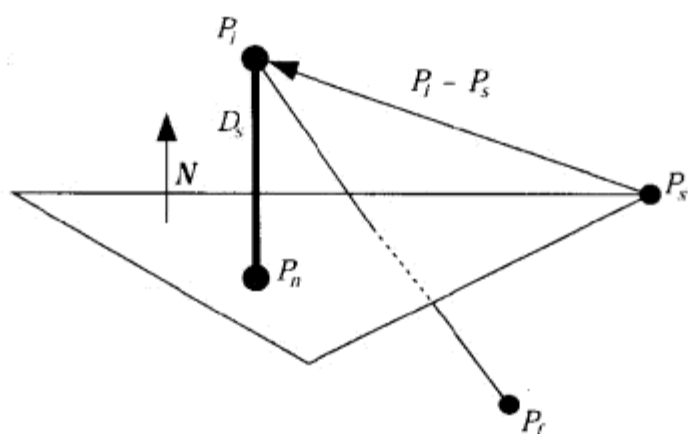


图 2.2.1 确定高度

面上最近的点

知道到平面的距离后，只需一步就可以获得平面上离起点  $P_i$  最近的点  $P_n$ ，如图 2.2.1 所示。我们知道，该点离起点的距离为  $D_s$ ，而该距离向量与面法线  $N$  平行。因此，可以这样来找到  $P_n$  点：

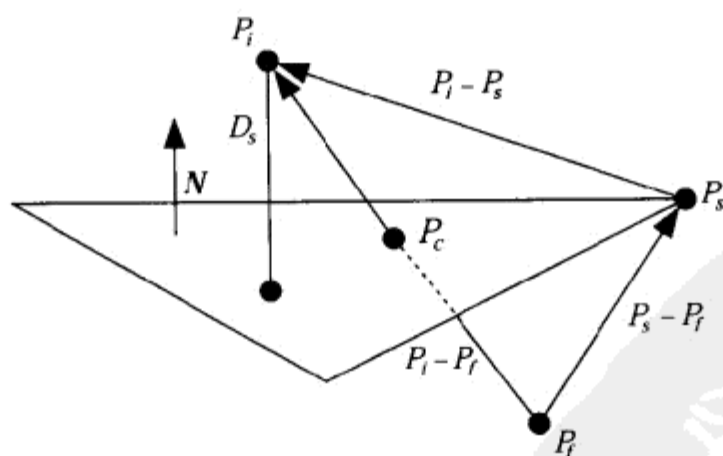
$$P_n = P_i - D_s N$$

法线向量的方向与我们距离向量相反，因此需要将起点减去距离向量。

## 2.2.2 找出碰撞点

当向量上的两点分别位于检测面的两边时，该向量必然会与平面相交。找出交点  $P_c$  便可以知道向量在什么位置穿过平面。如果碰撞检测代码没有指出交点，则可以按下面的方法来找。

您知道，碰撞点肯定位于向量所在的直线上。预先知道会发生碰撞可简化操作，因为我们知道该方程有一个解。如果检测射线与平面平行，则无法计算比例，因为这将导致被零除。我们可以利用查找碰撞点  $P_c$  时计算  $D_s$  的方法。图 2.2.2 说明了这种计算所需的信息。

图 2.2.2 查找碰撞点  $P_c$



由于我们知道碰撞点位于两点之间，因此可以通过计算它离  $P_i$  的距离来找到它。计算比例的公式如下：

$$R = ((P_i - P_s) \cdot N) / ((P_i - P_f) \cdot N)$$

使用前面计算得到的  $D_s$ ，上述公式变为：

$$R = D_s / ((P_i - P_f) \cdot N)$$

这两条线段相对于面法线的方向相同，这确保计算得到的比例不为负。有了比例后，就可以将其同  $P_i$  到  $P_f$  的向量相乘，得到碰撞点离  $P_i$  的距离。对于  $R=1$  的情况，无需做这样的计算，因为结果为  $P_f$ ；当  $R=0$  时，结果为  $P_i$ ；对于其他情况，可以使用下面的公式。

$$P_c = P_i + R(P_f - P_i)$$

### 2.2.3 到碰撞点的距离

虽然与  $D_s$  类似，但这并非到碰撞面的距离，因为它是沿行进路径的距离，而不是沿面法线方向的距离。如果行进路径离碰撞面很近，但方向几乎与碰撞面平行，则与高度相比，离碰撞点的距离将非常大。

计算飞行器的高度时，您将使用这样的计算方式，因为您不能保证面法线的方向。进行碰撞检测时，您将传递当前位置以及一个向下的非常大的向量（它足够长，能够与地面相交）；在计算与测试向量几乎垂直于小型多边形同碰撞面的交线时，也应这样做。在上述情况下，按全面介绍的方法计算得到的相对于碰撞面的高度  $D_s$  将非常小。

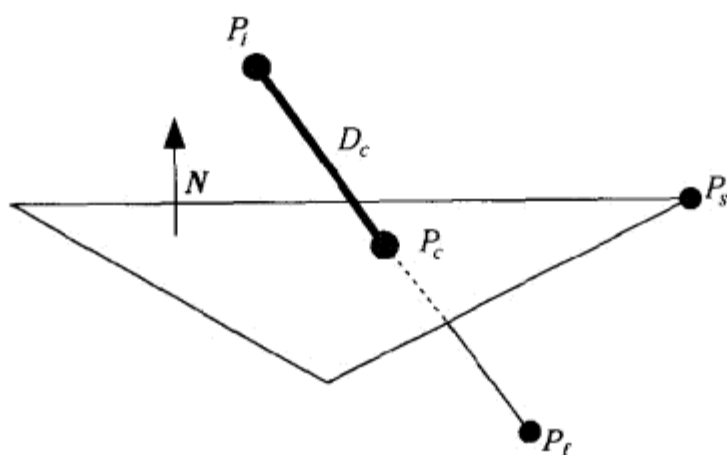


图 2.2.3 计算到碰撞点的距离

找到碰撞点后，计算碰撞点到起点的距离将非常容易，只需使用欧几里德方程即可。图 2.2.3 说明了所需的信息。到碰撞点的距离  $D_c$  为起点  $P_i$  到碰撞点  $P_c$  的向量的长度，而  $P_c$  已经在前面计算得到了。

$$D_c = |P_c - P_i|$$

另一种计算向量长度的方式是，它是向量的每个分量的平方和的平方根。大多数向量库





撞。您将需要不断重复这一过程，直到不再发生碰撞为止。每次发生碰撞后，反弹向量将变小（只要反弹不是在两个平行平面间进行的话）。

多次碰撞连续发生时，这是需要注意的衰减情况之一。如果碰撞正好发生在两个平面的交线上，则在进行第二次检测时，必须将面边缘上的点作为起点。如果您知道这种问题并有所考虑（假设正好位于面或边缘的任何东西都还未发生碰撞，并被视作位于平面的上方），则处理起来很容易。碰撞时，穿越距离总是大于零。

完成最后一次反射后，便可以计算新的速度向量，方法是将最后的碰撞点到最后的反射终点的向量规格化，然后将其与最初的速度向量相乘，如下所示：

$$V = \frac{(P_r - P_c) |P_i - P_f|}{|P_r - P_c|}$$

### 返回式碰撞 (Kickback collision)

有时候，您不希望玩家从碰撞面反弹出去，而希望它沿原路返回，如图 2.2.5 所示。对于这种情况，只要找出碰撞点后，计算将很简单。由于发生这种碰撞后，速度将保持不变，因此它也是完全弹性的。

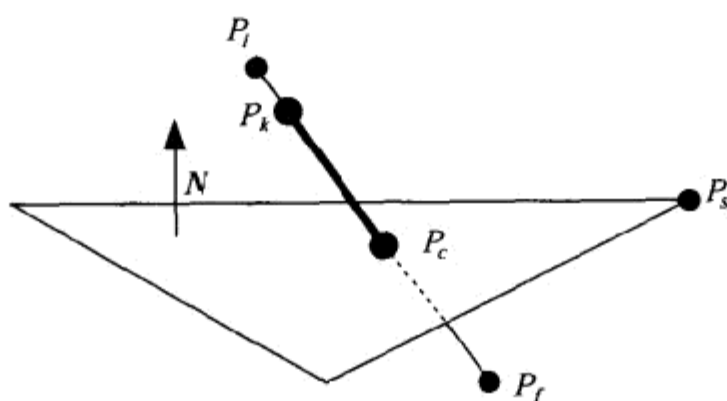


图 2.2.5 计算返回向量

返回终点  $P_k$  的计算方法是，计算终点  $P_f$  到碰撞点  $P_c$  的向量，然后将其与碰撞点相加。

$$P_k = P_c + (P_c - P_f)$$

合并同类项后，上述公式变为：

$$P_k = 2P_c - P_f$$

如果到达返回终点之前将发生其他碰撞，则返回式碰撞将出现反弹式碰撞中的问题。然而，在有些情况下，可以更早地脱离返回碰撞循环。如果碰撞点到终点的距离小于到起点的距离，则反弹终点将位于已经对碰撞情况进行过检测的区域，因此，没有必要再进行碰撞检测。

### 2.2.5 阻尼碰撞

如果碰撞是有摩擦或阻尼的，则处理向量是一定要小心。您将需要一个用作乘数的标量  $S$ ，它用于改变每次碰撞后的速度。 $S$  的值通常为 0~1。0 表示发生碰撞后，物体将停止；而 1 表示碰撞是完全弹性的，碰撞后物体的速度保持不变。通过将该标量增加到大于 1，可以向系统注入能量。返回式碰撞相当于标量  $S$  的值为负。

要正确地处理非弹性碰撞，必须取碰撞点  $P_c$  到反射终点  $P_r$  的向量的一部分（如图 2.2.6 所示），因为只在这一部分中由于碰撞而速度降低。下面的公式可用于计算阻尼碰撞的反射终点，该公式依赖于前面介绍的公式。

$$P_{slow} = P_c + (P_r - P_c)S$$

与此相对应，需要将其他单独存储的速度向量乘以同样的标量值，否则物体将在下一帧中恢复到原来的速度。如果像前面讨论的那样，碰撞导致发生另一次碰撞，则每次碰撞时都应使用该标量因子，以便在同一帧中模拟多次反弹的阻尼效果。

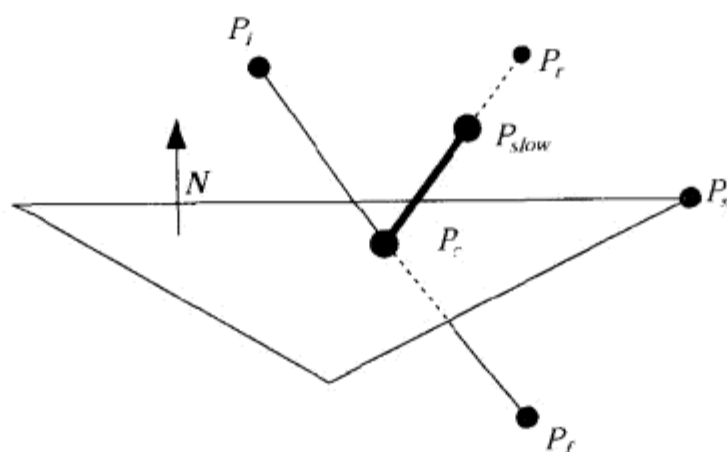


图 2.2.6 计算阻尼反射向量

#### 1. 线插值和面插值

线和面的一个有趣的地方是，在任何一组点之间进行加权平均插值得到的点都位于这组点定义的空间内。例如，对于直线，可以给它的两个端点指定权重，且两个端点的权重和为 1，则所有可能的加权平均插值点都将位于这两个点定义的直线上。添加一个点后，可形成一个平面，这将扩展上述规则。在此情况下，如果 3 个点的权重和为 1，则通过加权平均插值得到的点将位于这 3 个点定义的平面上。在平面的基础上可以添加其他的点，同时如果需要位于  $n$  维平面上的点，则可以增加维数。

这解释了为何前面的几个公式可以交替使用  $P_s$  和  $P_c$  的原因。这两个点都满足平面方程的要求。

另外，有趣的是权重无需位于 0~1 之间，而只需和为 1 即可。这使得加权平均插值得到的点位于点定义的线段或多边形的外面，但位于扩展后的线段或平面上。

## 2. 球与平面的碰撞

球与平面的碰撞比点与平面的碰撞要复杂些。对于这种碰撞的处理方式之一是使用比例。如果在起点  $P_i$  和基于线的碰撞点  $P_c$  之间划一条线，则当球与平面发生碰撞时，球心将位于这条线上。

当球刚接触到平面时，我们可以将  $P_i$  到  $P_c$  的直线同  $P_i$  到  $P_n$  的直线进行比较，以获得所需的信息。如果将第一条线投影到第二条线，则在  $P_i$  到  $P_n$  的直线上，球心离平面的距离刚好为球的半径。由于这条直线的长度为  $D_s$ ，因此可以计算出一个比例，该比例指出了球沿这条直线走了多远。这与我们使用比例来找到碰撞点  $P_c$  的情况类似。对于  $P_i$  到  $P_c$  的直线，比例与此相同，所以：

$$\frac{(P_c - P_b)}{(P_c - P_i)} = \frac{r}{D_s}$$

根据上述公式，可计算得到球心的位置，如下所示：

$$P_b = P_c - \frac{(P_c - P_i)r}{D_s}$$

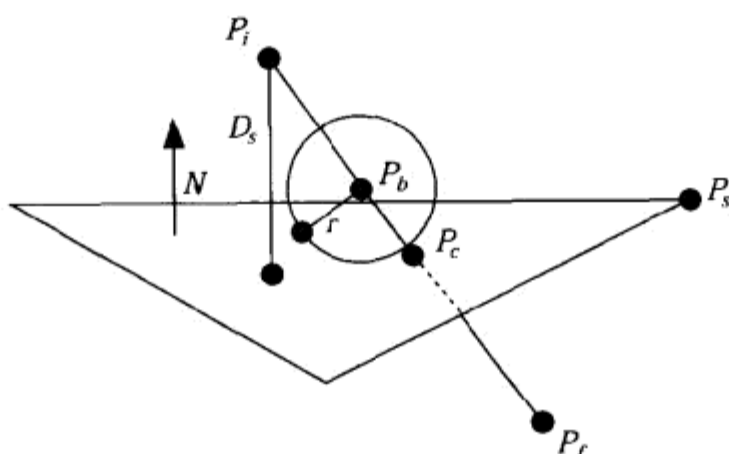


图 2.2.7 球与面的碰撞

图 2.2.7 图示说明了球相对于向量的位置。注意，球接触到平面时，并不会到达  $P_c$ ，除非  $P_i$  到  $P_c$  的直线垂直于平面。该向量与平面之间的角度越小，当球接触到平面时，它离  $P_c$  的距离将越大。

## 2.3 一种快速、健壮的计算 3D 线段交点的方法

Graham Rhodes, Applied Research Associates

grhodes@sed.ara.com

在游戏开发中，经常会遇到需要计算线段交点的情况。例如，对于简单的碰撞检测，计算线段交点很有帮助。假设两个三维空间的物体在不断移动。在一个时间步或动画帧中，每个物体都将沿直线从一个点移到另一个点。检测在时间步中两个物体是否会发生碰撞的最简单的方式是，看这两个线性路径之间的距离，如果它们的距离足够短（即小于两个物体的外接球的半径之和），则会发生碰撞。线段交点的其他常见用途有导航和移动规划（例如同 AI 系统结合使用时）、分色叠置片地图（map overlay）的创建以及地形/能见度估算等。

本文介绍一种健壮的、计算三维直线或线段交点（如果存在交点的话）的闭式（closed form）方法。如果没有交点，则计算两条线段上距离最近的点以及这两个点之间的向量。

### 2.3.1 这种算法健壮的原因

这里介绍的算法之所以健壮，原因有两个：首先，它没有任何特殊要求（如线段必须共面）；其次，它要求检查的例外情况很少。这种基本算法只要求检查两种特例，而这是数学上要求的，而不是推断所要求的。

### 2.3.2 问题描述

假设有两条三维线段，一条位于点  $\vec{A}_1 = [A_{1x} \ A_{1y} \ A_{1z}]^T$  和  $\vec{A}_2 = [A_{2x} \ A_{2y} \ A_{2z}]^T$  之间，而另条位于点  $\vec{B}_1 = [B_{1x} \ B_{1y} \ B_{1z}]^T$  和  $\vec{B}_2 = [B_{2x} \ B_{2y} \ B_{2z}]^T$  之间。我们的目标是找到这两条线段的交点  $\vec{P} = [P_x \ P_y \ P_z]^T$ （如果有的话）。如果没有交点，找到这条线段之间距离最近的点。图 2.3.1 对此做了说明。

距离最近的点分别为 C 和 D，它们可用于计算两条线段之间的最短距离。本文重点介绍如何找出距离最近的点，如果存在交点，则距离最近的点就是交点。

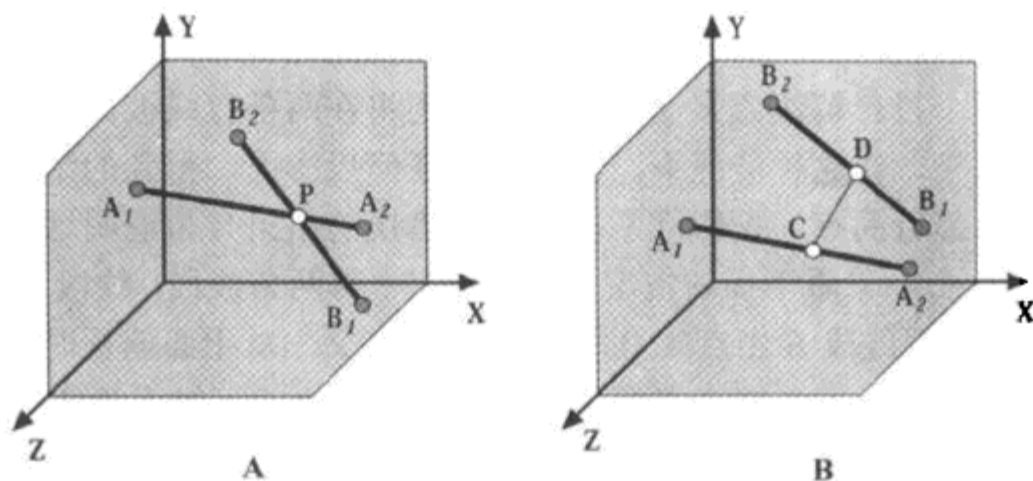


图 2.3.1 两条二维线段。A) 有交点；B) 没有交点

### 1. 分析

解决线段交点问题之前，进行一些分析将有所帮助。要正确地解决这种问题，我们面临的难题是什么？

对于任意一条三维直线，它可能与一个平面相交（如果直线不与平面平行的话）；但它很可能不会与另一条直线相交（即使这两条直线不平行，也不一定会相交）。基于上述分析可知，设计用于查找交点的算法将不是健壮的（即对于任何一对直线或线段，都能找出结果），因为大多数时候，这种算法都将失败。在实时的 3D 应用（如游戏）中，为获得一种健壮的算法，必须查找两条直线之间距离最近的点。

任何学习过基本平面几何课程的学生都求过二维直线的交点，因此探讨一下求二维直线交点和求三维直线交点之间的关系将会有所帮助。在二维空间中，不平行的两条直线都将相交于一点。为说明三维空间的情况，假设某个平面包含定义直线  $A$  的两个点和定义直线  $B$  的第一个点。直线  $A$  位于该平面中，定义直线  $B$  的第一个点也是如此。因此，这两条直线的交点肯定位于该平面上，因为交点位于直线  $A$  上。同时，交点还位于直线  $B$  上，因此直线  $B$  有两个点位于该平面上。由于直线  $B$  有两个点位于该平面上，因此直线  $B$  位于该平面中。

这里的重要结论是，如果两条直线之间有交点，则这两条直线必然共面。因此，如果两条三维直线之间有交点，则求三维直线交点的问题与求二维直线交点的问题相同。

### 2. 幼稚的算法

对于交点问题，一种幼稚的算法（这种算法有问题）是，将两条直线投影到一个标准坐标平面（ $XY$ 、 $YZ$  或  $XZ$ ）上，然后在平面中解决这种问题。就实现方面而言，这种方法的主要困难在于要选择一个合适的平面，将直线投影到该平面中。如果其中任何一条线段都不与任何坐标平面平行，则可在任何坐标平面内解决这个问题。然而，如果至少有一条线段与坐标平面平行，则要求大量的逻辑，这是无法接受的。这种方法的一个变体（它不那么幼稚，但仍然存在问题）是，根据定义直线的 4 个点中的 3 个定义一个平面，并将这 4 个点都投影到该平面中，然后在平面中解决问题。在直线之间确实有交点时（这样的情况不常见），这种方法确实能得到正确的结果。

这两种通过二维投影来解决交点问题的方法缺乏的重要功能之一是，不能指出是否存在



三维交点。另外，也没有指出距离最近的点。必须通过反推来获得这种重要的信息。



这两个投影解决方案变种的最大问题将在两条直线彼此擦肩而过但没有相交时显现出来。在这种情况下，无论使用任何投影面，得到的结果都不一定是正确的、距离最近的点。投影常常得到完全错误的结果！为说明这一点（用语言说明起来很困难），来看下面的情况。三维空间中有两条线段，线段  $A$  由点  $(0, 0, 0)$  和  $(1, 0, 0)$  定义，而线段  $B$  由点  $(1, 0, 1)$  和  $(1, 1, 1)$  定义。从上面俯视时（相当于将线段投影到  $XY$  平面上），二维交点为  $(1, 0, 0)$ ，距离最短的三维点为  $(1, 0, 0)$  和  $(1, 0, 1)$ 。对于这个问题，这里的答案是正确的。然而，如果从其他角度看，二维交点将可能出现在两条线段的任何位置。通过将投影到二维空间中得到的答案，反推到三维直线中得到的距离最短的对有无数个，显然这是不正确的。附带光盘中的测试代码有助于理解这里存在的问题，它让您能够旋转视图，从不同的视角来查看两条线段，同时显示距离最近的点，让您能够将其同观察区看到的交点进行比较。

下一节将推导出一种计算点  $C$  和  $D$  的闭式（closed form）解决方案，该解决方案没有对线段的位置做任何假设。该解决方案确实处理了两种特殊情况，但即使在其他替代方法中，这些情况也是不可避免的。

### 2.3.3 推导闭式解决方案方程

#### 1. 计算两条直线之间距离最近的点

在三维空间中，可将直线方程视为单个标量值（参数）的向量函数。为推导计算两条三维直线之间距离最近的两个点的解决方案，我们首先写出第一条直线上任意一个点  $\vec{C} = [c_x \ c_y \ c_z]^T$  的方程，如方程 2.3.1 所示。

$$\vec{C} = \vec{A}_1 + s\vec{L}_A, \text{ 其中 } \vec{L}_A = (\vec{A}_2 - \vec{A}_1) \quad (2.3.1)$$

方程 2.3.1 的基本含义是，第一条直线上任意一个点的坐标等于定义该直线的第一点的坐标加上第一个点到第二个点的向量与标量参数  $s$  的乘积。如果  $s$  为 0，则为第一个点；如果  $s$  等于 1，则为第二个点。对于第二条直线上的点  $\vec{D} = [d_x \ d_y \ d_z]^T$ ，可以书写类似的方程，如方程 2.3.2 所示。

$$\vec{D} = \vec{B}_1 + t\vec{L}_B, \text{ 其中 } \vec{L}_B = (\vec{B}_2 - \vec{B}_1) \quad (2.3.2)$$

在上述方程中， $t$  是另一个标量参数，其物理含义与  $s$  相同。如果参数  $s$  和  $t$  可以是任意的，则使用它们来计算  $\vec{C}$  和  $\vec{D}$  时，得到的将是直线而不是线段。对于线段，参数  $s$  和  $t$  的取值范围为  $0 \leq s, t \leq 1$ 。现在假设  $s$  和  $t$  的值可以是任意的，后面将考虑线段的情况。

如果可以找到这样的  $s$  和  $t$  值，即使得  $\vec{C}$  和  $\vec{D}$  重叠，则说明两条三维直线相交。通常情况下，相交的情况比较罕见，然而我们需要一种计算  $\vec{C}$  和  $\vec{D}$  为距离最近的点时  $s$  和  $t$  的值的



方法。接下来的推导演示如何得到这样的  $s$  和  $t$  值。

首先，将方程 2.3.2 和 2.3.1 相减，得到表示从点  $\vec{C}$  到  $\vec{D}$  的向量的方程：

$$\vec{C} - \vec{D} = -\vec{AB} + s\vec{L}_A - t\vec{L}_B = [0 \ 0 \ 0]^T, \text{ 其中 } \vec{AB} = \vec{B}_1 - \vec{A}_1 \quad (2.3.3)$$

由于我们希望点  $\vec{C}$  和  $\vec{D}$  重叠，因此将这两点之间的向量设置为零向量。方程 2.3.3 右边可表示为如下矩阵方程：

$$\begin{bmatrix} L_{Ax} & -L_{Bx} \\ L_{Ay} & -L_{By} \\ L_{Az} & -L_{Bz} \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} AB_x \\ AB_y \\ AB_z \end{bmatrix} \quad (2.3.4)$$

方程 2.3.4 有三行，每行对应一个坐标方向，但有两个未知数——标量值  $s$  和  $t$ 。这是一个典型的由多种因素决定（受约束）的系统。仅当右边的系数矩阵的秩为 2 时，该方程才有精确解，在这种情况下，三个方程将等价于两个独立的方程，因此  $s$  和  $t$  将有精确解。从几何上说，有精确解意味着两条直线相交，且共面。因此，两条共面的三维直线只有一个交点。

方程 2.3.4 两边的差为  $\vec{C}$  和  $\vec{D}$  之间的向量，它也是对任意  $s$  和  $t$  值，方程 2.3.4 的误差向量。我们通过设置  $s$  和  $t$  的值，使该向量的长度最短来获得距离最近的点。

使得  $\vec{C}$  和  $\vec{D}$  之间的距离最短的  $s$  和  $t$  值正是方程 2.3.4 的最小线性平方解。从几何上说，最小二乘方解得到的是点  $\vec{C}$  和  $\vec{D}$ 。当两条直线共面但不平行时，该算法得到的将是实际的交点。方程 2.3.4 可以写成下面的形式：

$$M\vec{s} = \vec{b}, \text{ 其中 } \vec{s} = [s \ t]^T \quad (2.3.5)$$

对于多因素决定的系统，找出最小二乘方解的方法之一是解标准方程而不是原来的系统 [Golub96]。标准方程方法适合于这个问题，但对于涉及线性方程系统的一般性问题，则可能存在疑问。我们将左边和右边同时左乘系数矩阵  $M$  的转置矩阵来获得标准方程。对于我们这里的问题，标准方程如方程 2.3.6 所示。

$$M^T M \vec{s} = M^T \vec{b}, \text{ 其中 } M^T \text{ 为 } M \text{ 的转置矩阵} \quad (2.3.6)$$

方程 2.3.6 具备所需的特性，即将系统简化为两个方程组成的系统的解，而要用代数方法获得  $s$  和  $t$  的解正好需要两个方程。我们来完成推导方程 2.3.4 的标准方程的过程。将方程 2.3.6 扩展后，得到如下标准方程：

$$\begin{bmatrix} L_{Ax} & L_{Ay} & L_{Az} \\ -L_{Bx} & -L_{By} & -L_{Bz} \end{bmatrix} \begin{bmatrix} L_{Ax} & -L_{Bx} \\ L_{Ay} & -L_{By} \\ L_{Az} & -L_{Bz} \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} L_{Ax} & L_{Ay} & L_{Az} \\ -L_{Bx} & -L_{By} & -L_{Bz} \end{bmatrix} \begin{bmatrix} AB_x \\ AB_y \\ AB_z \end{bmatrix} \quad (2.3.7)$$

代入矩阵代数：

$$\begin{bmatrix} \vec{L}_A \cdot \vec{L}_A & -\vec{L}_A \cdot \vec{L}_B \\ -\vec{L}_A \cdot \vec{L}_B & \vec{L}_B \cdot \vec{L}_B \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} \vec{L}_A \cdot \vec{AB} \\ -\vec{L}_B \cdot \vec{AB} \end{bmatrix} \quad (2.3.8)$$

或者，通过定义一系列新的标量变量，可简化为：

$$\begin{bmatrix} L_{11} & L_{12} \\ L_{12} & L_{22} \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} r_A \\ r_B \end{bmatrix} \quad (2.3.9)$$

这是一个简单的  $2 \times 2$  矩阵。本节余下的内容将以代数方式解上述方程，以获得  $s$  和  $t$  的封闭解。对于方程 2.3.9，解的方法有很多，其中包括克莱姆法则[O'Neil87]和高斯消元法[Golub96]。从理论上说，克莱姆法则很有趣，但代价很高，需要执行大约  $(n+1)!$  次乘法和除法运算。高斯消元法的代价低些，需要执行  $n^3/3$  次乘除运算。对于线性方程系统，还有其他的解法，当系统非常大时，这些方法更可靠，通常速度也更快。其中包括诸如 QR 因式分解（用于中等规模的矩阵）等高级直接解法以及迭代法（用于大型稀疏矩阵）。我将使用高斯消元法来推导解，对于  $2 \times 2$  的矩阵来说，这种方法的代价比克莱姆法则低些。下面我们执行行消除，得到一个上三角矩阵。行消除步骤如下。对于方程 2.3.9，将第 2 行修改为第 2 行与第 1 行的  $L_{12}/L_{11}$  倍的差，得到方程 2.3.10。

$$\begin{bmatrix} L_{11} & L_{12} \\ L_{12} - L_{11} \frac{L_{12}}{L_{11}} & L_{22} - L_{12} \frac{L_{12}}{L_{11}} \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} r_A \\ r_B - r_A \frac{L_{12}}{L_{11}} \end{bmatrix} \quad (2.3.10)$$

对方程 2.3.10 进行化简，并将第 2 行乘以  $L_{11}$ ，得到如方程 2.3.11 所示的上三角矩阵：

$$\begin{bmatrix} L_{11} & L_{12} \\ 0 & L_{11}L_{22} - L_{12}^2 \end{bmatrix} \begin{bmatrix} s \\ t \end{bmatrix} = \begin{bmatrix} r_A \\ L_{11}r_B - L_{12}r_A \end{bmatrix} \quad (2.3.11)$$

根据方程 2.3.11，可以得到  $t$  和  $s$  的解：

$$t = \frac{L_{11}r_B - L_{12}r_A}{L_{11}L_{22} - L_{12}^2} \quad (2.3.12)$$

$$s = \frac{r_A - L_{12}t}{L_{11}} \quad (2.3.13)$$

需要指出的是，在某些退化情况下，方程 2.3.12 和 2.3.13 无效，这要求我们以有限的方式使用允许容错。如果线段 A 的长度为 0，则方程 2.3.13 无效；如果任何一条线段的长度为 0 或两条线段平行，则方程 2.3.12 无效。在这些情况下，将导致被零除的异常。后面的“特殊情况”一节对此做了更详细的讨论。

对于  $2 \times 2$  的矩阵，就计算成本而言，使用高斯消元法和克莱姆法则来求解  $s$  和  $t$  之间的差别在于，使用高斯消元法时，求解  $s$  的值需要执行 1 次乘法、1 次除法和 1 次减法运算；而使用克莱姆法则时，需要执行 4 次乘法、1 次除法和 2 次减法运算。

现对推导过程总结如下。对于有点  $\bar{A}_1$  和  $\bar{A}_2$  定义的直线 A 以及由点  $\bar{B}_1$  和  $\bar{B}_2$  定义的直线 B，定义下述中间变量：

$$\bar{L}_A = (\bar{A}_2 - \bar{A}_1); \quad \bar{L}_B = (\bar{B}_2 - \bar{B}_1); \quad \bar{AB} = \bar{B}_1 - \bar{A}_1 \quad (2.3.14)$$

以及：

$$\begin{aligned} L_{11} &= \vec{L}_A \cdot \vec{L}_A; & L_{22} &= \vec{L}_B \cdot \vec{L}_B; & L_{12} &= -\vec{L}_A \cdot \vec{L}_B \\ r_A &= \vec{L}_A \cdot \vec{AB}; & r_B &= -\vec{L}_B \cdot \vec{AB} \end{aligned} \quad (2.3.15)$$

计算参数  $s$  和  $t$  的值，这两个参数定义了距离最近的两个点：

$$t = \frac{L_{11}r_B - L_{12}r_A}{L_{11}L_{22} - L_{12}^2} \quad (2.3.16)$$

和

$$s = \frac{r_A - L_{12}t}{L_{11}} \quad (2.3.17)$$

第一条直线上离第二条直线最近的点为：

$$\vec{C} = \vec{A}_1 + s\vec{L}_A \quad (2.3.18)$$

而第二条直线上离第一条直线最近的点为：

$$\vec{D} = \vec{B}_1 + t\vec{L}_B \quad (2.3.19)$$

可以这样认为，位于两个距离最近的点中间点是离两条直线/线段最近的点：

$$\vec{P} = (\vec{C} + \vec{D})/2 \quad (2.3.20)$$

当然，如果两条直线相交，则  $\vec{P}$  将为交点。

## 2. 特殊情况

就两条三维线段之间距离最近的点而言，只有两个可能的特殊情况。第一种情况是至少一条线段退化——由两个重叠在一起的点定义，即  $\vec{A}_1$  和  $\vec{A}_2$  重叠或  $\vec{B}_1$  和  $\vec{B}_2$  重叠。我们将这种情况称为线段退化特例。第二种情况是两条线段平行，这被称为线段平行特例。

将线段退化特例同前面推导出的方程关联起来很容易。注意，方程 2.3.15 中定义的变量  $L_{11}$  等于线段  $A$  长度的平方；而变量  $L_{22}$  等于线段  $B$  长度的平方。如果这两个变量中的任何一个为 0，则说明相应的线段已退化，这样，方程 2.3.9 中的矩阵的行列式的结果为 0，无法找到  $s$  和  $t$  的解。当  $L_{11}$  或  $L_{22}$  为 0 时， $L_{12}$  也为 0。

判断线段  $A$  是否退化的标准方式之一是：

```
bool line_is_degenerate = L11 < ε2 ? true : false;
```

其中  $\epsilon$  是一个很小的数，如  $10^{-6}$ 。将  $\epsilon$  设置为一个诸如  $10^{-6}$  这样的值比将其设置为小得多的计算机小正数更合适。

当线段  $A$  和  $B$  都已退化时，则可以将点  $\vec{C}$  设置为点  $\vec{A}_1$ ，将点  $\vec{D}$  设置为点  $\vec{B}_1$ 。当只有线段  $A$  退化时，可将点  $\vec{C}$  设置为点  $\vec{A}_1$ ，并计算线段  $B$  上离点  $\vec{C}$  最近的点，这点就是点  $\vec{D}$ 。在这种情况下，只需根据方程 2.3.21 计算参数  $t$  即可。

$$-\vec{L}_B t = \vec{AB} \quad (2.3.21)$$

方程 2.3.21 是线段 A 退化时, 方程 2.3.4 的简化结果。该方程要求我们找到最小二乘方解。这里无需推导, 只需使用标准方程即可得到如下最小二乘方解:

$$t = \frac{-\vec{L}_B \cdot \vec{AB}}{\vec{L}_B \cdot \vec{L}_B} = \frac{r_B}{L_{22}} \quad (2.3.22)$$

点  $\vec{D}$  可使用方程 2.3.2 来计算得到。

当只有线段 B 退化时, 将点  $\vec{D}$  设置为点  $\vec{B}_1$ , 并计算线段 A 上离点  $\vec{D}$  最近的点, 这点就是点  $\vec{C}$ 。在这种情况下, 只需根据方程 2.3.23 计算参数  $t$  即可, 该方程类似于方程 2.3.21。

$$\vec{L}_A s = \vec{AB} \quad (2.3.23)$$

$s$  的解为:

$$s = \frac{r_A}{L_{11}} \quad (2.3.24)$$

注意, 方程 2.3.24 是  $t$  等于 0 时方程 2.3.13 的结果。由于  $t$  等于 0, 因此这里的推导结果与线段没有退化的情况一致。

当然, 对于只有一条线段退化的情况, 一种不错的方法是编写一个子程序, 并将它用于只有线段 A 退化和只有线段 B 退化的情况。只要正确地处理变量, 便可以使用方程 2.3.22 或方程 2.3.24 来完成这样的工作。在附带光盘提供的实现中, 使用的是方程 2.3.24, 它通过传递参数使得退化的线段总被视为线段 B, 而未退化的线段总被视为线段 A。



虽然没有线段退化的情况那么明显, 但将线段平行特例与前面推导出的方程关联起来也很容易。 $L_{12}$  为  $-\vec{L}_A \cdot \vec{L}_B$ , 当线段平行时, 上述点积将为  $\vec{L}_A$  长度和  $\vec{L}_B$  长度的乘积再取负。方程 2.3.9 中的矩阵对应的行列式的结果为  $L_{11}L_{22} - L_{12}^2$ , 当  $L_{12}$  的大小等于  $\vec{L}_A$  长度和  $\vec{L}_B$  长度的乘积时, 该结果将为 0。因此, 当线段平行时, 方程 2.3.9 是奇异的, 无法获得  $s$  和  $t$  的解。

对于两条平行的直线, 直线 A 上的所有点到直线 B 的距离都相等。如果必须计算直线 A 和直线 B 之间的距离, 只需将点  $\vec{C}$  设置为  $\vec{A}_1$ , 然后使用方程 2.3.22 和 2.3.2 找到点  $\vec{D}$  即可。这样, 点  $\vec{C}$  和  $\vec{D}$  之间的距离就是两条直线之间的距离。对于线段的情况, 将在下一节介绍如何处理。

### 3. 编码效率

为提高编码效率, 应首先检查直线退化的情况, 然后检查直线平行的情况。这样, 如果至少有一条直线退化, 则无需计算方程 2.3.14 和 2.3.15 中一些出于方便的目的而设置的变量。

## 2.3.4 线段

前面两节处理的是直线，然而在游戏开发中需要处理线段的情况更常见。那么，如果对前面的推导结果进行调整，以处理线段的情况呢？

## 1. 线段不平行

如果根据方程 2.3.12 和 2.3.13 得到的  $s$  和  $t$  的值都位于范围  $[0, 1]$  内，则无需做任何其他的工作，因为线段的结果与直线相同。如果  $s$  和  $t$  的值至少有一个不在上述范围内，则必须对结果进行调整。线段不平行时，存在两种可能性： $s$  或  $t$  位于范围  $[0, 1]$  之外； $s$  和  $t$  都位于范围  $[0, 1]$  之外。图 2.3.2 说明了这两种情况。

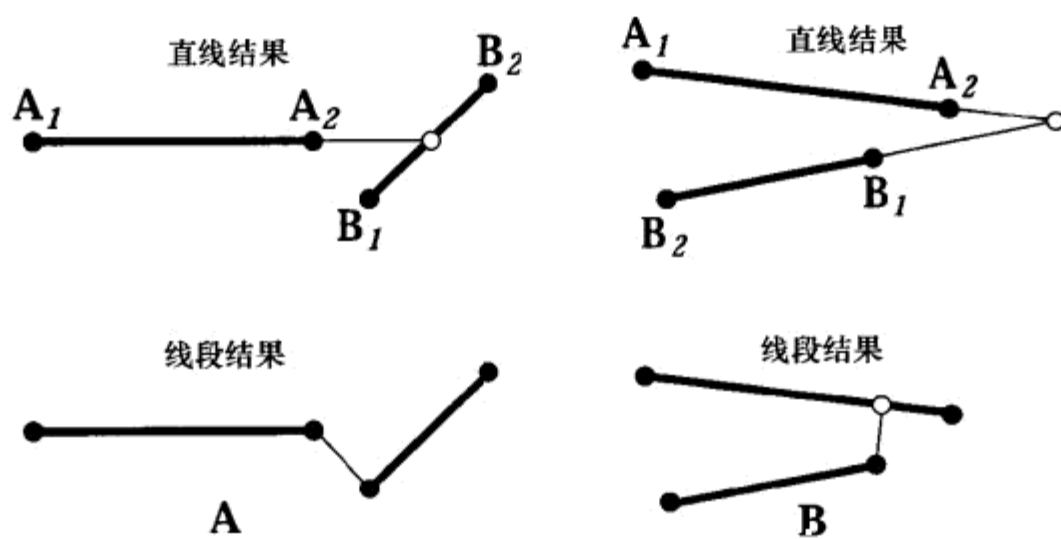


图 2.3.2 线段。A)  $s$  和  $t$  中的一个位于范围  $[0, 1]$  之外；B)  $s$  和  $t$  都位于范围  $[0, 1]$  之外

对于  $s$  或  $t$  位于范围  $[0, 1]$  之外的情况（如图 2.3.2a 所示），需要执行下述操作：

- (1) 将超出范围的参数定位到  $[0, 1]$ ；
- (2) 计算新参数对应的点，这点离另一条线段的距离最近；
- (3) 按计算线段间距离最近点的方法，找出另一条线段上离上述点距离最近的点。这便是另一条线段上距离最近的点。

在最后一步，只需将方程 2.3.22 的结果定位到  $[0, 1]$  内，然后计算相应的点即可。

$s$  和  $t$  都不在范围  $[0, 1]$  的情况（如图 2.3.2B 所示）稍微要复杂些，但处理方式与前面相同，只是需要决定先处理哪条线段。例如，对于图 2.3.2B 所示的情况，如果选择先处理线段 A，则第 2 步的结果为点  $\bar{A}_2$ 。然后第 3 步得到的结果将为点  $\bar{B}_1$ ——线段 B 上离点  $\bar{A}_2$  最近的点。

显然， $\bar{A}_2$  和  $\bar{B}_1$  并非正确的  $\bar{C}$  和  $\bar{D}$ 。对于  $\bar{D}$  而言，将其设置为  $\bar{B}_1$  是正确的，但在线段 A 上，存在

在离线段 B 的距离比点  $\bar{A}_2$  更近的点。事实上，无论第 3 步计算的是  $\bar{C}$  还是  $\bar{D}$ ，得到的结果都

将是正确的。错误的结果来自第 2 步。因此，我们可以根据第 3 步的结果来计算另一个距离最近的点。这样，对于  $s$  和  $t$  都位于范围  $[0, 1]$  的情况，可以采取下面的步骤：

- (1) 选择一条线段，将其超出范围的参数定位到  $[0, 1]$ ；

(2) 计算新参数对应的点，这点并不一定离另一条线段的距离最近；

(3) 按计算线段间距离最近点的方法，找出另一条线段上离上述点最近的点，这是另一条线段上距离最近的点；

(4) 找出第一条线段上离第 3 步得到的点最近的点，这是第一条线段上距离最近的点。

对于图 2.3.2B 的情况，如果首先选择线段  $B$ ，则第 2 步的结果将为  $\bar{B}_1$ ，而第 3 步得到的结果将为  $\bar{A}_1$  和  $\bar{A}_2$  的点。在这种情况下，无需进行第 4 步。这里的实现没有检查这种情形。

## 2. 线段平行

线段平行时，存在两种可能的情况，图 2.3.3 说明了这两种情况。首先，可能存在惟一两个点，它们之间的距离最短，如图 2.3.3A 所示；当将这两条线段投影到与它们平行的直线上时，如果它们不重叠，则属于这种情况。其次，可能存在多对点，它们之间的距离是最近的，如图 2.3.3B 所示。在这种情况下，可以选择两条垂直灰线之间的任何一对距离最近的点。对于平行且重叠的线段，附带光盘提供的实现选择的是两条灰色正中间的一对距离最近的点，即每条线段的重叠部分的中点。

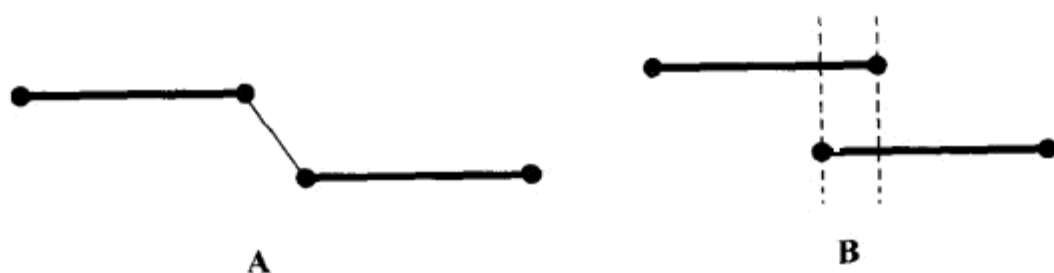


图 2.3.3 平行线段。A) 存在惟一的两个点；B) 存在很多对点

需要指出的是，当两条线段平行或几乎平行且稍微旋转时，使用这种算法计算得到的距离最近的点常常以无规律的方式移动。在这种情况下，该算法管用，但得到的结果常常是令人迷惑且质疑的，因为距离最近的点将在线段两端之间来回跳动。图 2.3.4 说明了这一点。

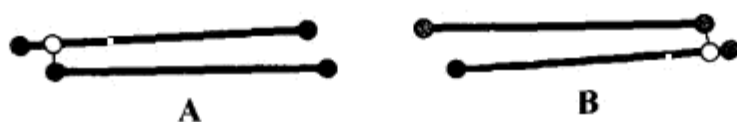


图 2.3.4 线段几乎平行时，距离最近的点漂移不定。

A) 距离最近的点位于左端；B) 距离最近的点位于右端







从图 2.3.4A 可知, 距离最近的点位于最左端, 但随着线段旋转到平行时, 距离最近的点将跳到重叠的部分。然后, 随着线段继续旋转, 距离最近的点将跳到最右端, 如图 2.3.4B 所示。在有些游戏应用程序中, 这样的行为可能引发问题。为处理这样的问题, 可以在线段几乎平行或完全平行时, 采用其他的方法来选择距离最近的点。例如, 可以实现一条这样的规则, 即当线段几乎平行时(如角度小于  $5^\circ$ ), 随机地选择  $\bar{A}_1$  附近的某个点作为线段  $A$  上距离最近的点。为避免角度为  $5^\circ$  时出现不稳定的情况, 需要在某个角度范围内(如  $5^\circ \sim 10^\circ$ ) 将该随机点同另一个计算得到的值进行混合—— $5^\circ$  时, 随机点的权重为 100%, 而  $10^\circ$  时, 其权重为 0%。这种解决方案将增加该算法的开销。当然, 还有其他的方法, 这些方法可能更简单、开销更低, 也更可靠。附带光盘中提供的实现没有对这样的行为进行控制。

### 2.3.5 实现描述

实现中包括 4 个 C 语言函数, 这些函数位于文件 `lineintersect_utils.h` 和 `lineintersect_utils.cpp` 中。主要的接口是函数 `IntersectLineSegments`, 它接受定义两条线段的参数, 并返回点  $\bar{C}$ 、 $\bar{D}$  和  $\bar{P}$  以及点  $\bar{C}$  到  $\bar{D}$  的向量。该函数还接受一个指出是否应将线段视为直线的参数以及一个用于检测退化特例和平行特例的参数容错参数。在该实现的外面, 可以使用从  $\bar{C}$  到  $\bar{D}$  的向量来计算直线之间的距离。需要指出的是, 对于线段, 该向量并不一定垂直于任何线段; 但对于直线且至少有一条未退化的情况, 则该向量必然垂直于未退化的直线。辅助函数如下:

- `FindNearestPointOnLineSegment`: 计算线段上离三维空间中给定的点距离最近的点;
- `FindNearestPointOnLineSegment`: 计算线段平行时表示点  $\bar{C}$  和  $\bar{D}$  的值(可能是惟一的, 也可能不是);
- `AdjustNearestPoints`: 对直线情况下得到的点  $\bar{C}$  和  $\bar{D}$  的值进行调整, 使之成为线段情况下的解。

代码中包含注释, 为参考正文提供了索引。



另外, 还提供了一个测试程序, 该程序名为 `line_intersection_demo`。对于 OpenGL, 必须将该演示程序链接到 GLUT 库。工程文件是使用 Microsoft Visual C++ 6.0 编写的。要将其移植到其他支持 OpenGL 和 GLUT 的系统中, 应不会太困难。

### 2.3.6 可优化的地方

编写实现的源代码时, 作者很小心, 但未针对特定的处理器或指令集进行优化。对于给定的平台, 任何实现都有很多可优化的地方。就这个实现而言, 优化空间最大的方面可能是向量化。在这些代码中, 很多运算都涉及到对向量的所有三个分量进行乘或加/减运算。这是向量化方面可优化的主要地方。另外, 如果指令集支持诸如点积等高级运算, 则在计算方程 2.3.15 时可充分利用这种特性。要最大限度地提高性能, 强烈建议您使用专业代码剖析



(profiling) 工具找出目标平台中的瓶颈和可优化的地方。



本文以及附带光盘中的实现是非常严密的, 考虑到了各种可能出现的情况。总体而言, 代码的效率很高, 但对于直线在线段外相交的情况(即  $s$  和  $t$  至少有一个不在范围  $[0, 1]$  内), 计算真正的距离最短的点的方法不一定是效率最高的。实际上, 本文解决的距离最近的点的问题是最为简单的, 通常情况下, 当在这个问题的基础上加上约束条件时, 计算的成本将增加。除处理器/平台特定的优化外, 还可以根据应用的实际需要, 删除该实现中多余的部分。例如, 如果无需处理线段的情况, 则可删除该实现中与线段相关的所有代码。而只需在找到的距离最近的点位于线段端点之间时, 让函数返回布尔值 `true`, 否则返回 `false`。

### 2.3.7 结论

---

本文讨论的算法是严密的, 能够处理任何直线交点问题。您可能需要对该算法进行调整, 这取决于您将如何使用交点, 如处理两条线段几乎平行时出现的异常情况, 或者当您只需要处理直线时, 删除对线段进行处理的代码。真诚地希望读者能够通过阅读这里关于直线/线段相交的讨论以及光盘中附带的代码受益。

### 2.3.8 参考文献

---

[Golub96] Golub, Gene H., and Charles F. van Loan, *Matrix Computations, Third Edition*, The Johns Hopkins University Press, 1996.

[O'Neil87] O'Neil, Peter V., *Advanced Engineering Mathematics, Second Edition*, Wadsworth Publishing Company, 1987.



## 2.4 反向弹道计算

Aaron Nicholls, Microsoft

aaron\_feedback@hotmail.com

在游戏开发中，一个常出现的问题是计算弹道。在最常见的情况下，导弹的速度和方向是已知的，需要计算导弹在给定时间的位置以及是否会与其他物体相撞。这是一个简单的迭代问题，但大部分游戏要求的并不仅仅是这些。很多情况下，还需要解决这个问题的逆问题，即在给定一系列常数（如重力加速度、起点、终点）的情况下，计算导弹的发射角度和速度。另外，有了这种问题的解决方案后，便可以将其用作解决更复杂的此类问题的框架。

本文假设读者熟悉基本的二维/三维变换、基本积分和三角函数。

### 1. 简化探讨的问题

简化问题的方式有多种，我们可以首先将问题从三维简化为二维。在导弹的初始速度和方向给定的情况下，如果作用力只有地球引力（通常可假定为常数），则弹道将为共面的抛物线。因此，通过将该共面的弹道变换到二维空间（ $x$  和  $y$ ）中，可以极大地简化问题。另外，通过将起点变换为原点，可以消除大多数方程中的初始  $x$  值和  $y$  值，重点将在终点的坐标上。图 2.4.1 是一个弹道范例，它被旋转到  $xy$  平面中，同时起点被平移到原点。

另外，我们需要确定要解决的到底是什么问题。就这里而言，变量为初始速度、仰角以及两点之间的  $x$  距离和  $y$  距离。在这 4 个变量中的 3 个已知（还有 1 个未知）时，我们的目标是推导一个方程，用 3 个已知的值计算出未知的值。

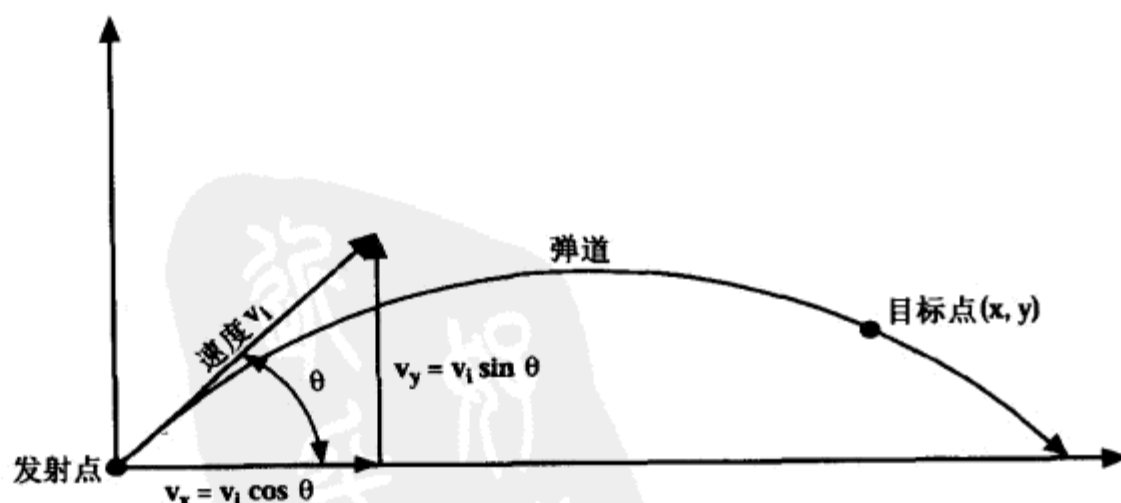


图 2.4.1 二维空间中两点之间的弹道

然而，存在多个未知量的情况很常见。在这样的情况下，最佳的解决方案是将其其中的一些变量的值设置为常量，从而消除该变量。例如，我们常常知道发射点和目标点的位置，需要提供初始速度和仰角。在这种情况下，可以将初始速度设置为最大的可能速度  $v_{max}$ ，从而消除这个变量。这样，只有一个量是未知的，因此只需使用  $v_i$ 、 $x$  和  $y$  来计算仰角  $\theta$  即可。本文后面的“处理多个变量”一节将更详细地讲述这种技巧及其使用指南。

## 2. 将位置和速度定义为时间的函数

将问题简化为二维后，可以确定每一维的速度和加速度。假设初始速度为  $v_i$ 、仰角为  $\theta$ 、重力加速度为  $g$ ，则沿  $x$  轴和  $y$  轴的初始速度如下：

$$v_{xi} = v_i \cos \theta$$

$$v_{yi} = v_i \sin \theta$$

由于导弹受到的作用力只有地球引力，因此可以假设水平速度 ( $v_x$ ) 保持不变，地球引力只影响垂直速度 ( $v_y$ )。垂直速度和水平速度可被表示为：

$$v_x = v_i \cos \theta \quad (2.4.1)$$

$$v_y = v_i \sin \theta - gt \quad (2.4.2)$$

接下来，将上述速度公式合并起来，以计算导弹在给定时间的位置（假设发射点为原点）。

$$x = \int v_i \cos \theta \, dt$$

$$\rightarrow x = v_i t \cos \theta \quad (2.4.3)$$

$$y = \int (v_i \sin \theta - gt) \, dt$$

$$\rightarrow y = v_i t \sin \theta - \frac{1}{2} g t^2 \quad (2.4.4)$$

### 2.4.1 一种特殊情况

#### 1. 发射点和目标点的海拔高度相同

处理这种问题的普通情况之前，来看一种更简单的情况，这将帮助我们理解如何解决更普遍的问题。一种常见的特例是，发射点和目标点的  $y$  坐标相同。一个这样的例子是，游戏中的地面是平坦的，地面上的大炮向地面目标开火。在这种情况下，发射点和目标点之间的垂直距离为 0，因此可以将  $y$  设置为 0，从而简化垂直位置方程。这让我们能够简化方程 2.4.4，从而求出时间  $t$ 、初始速度  $v_i$  或仰角  $\theta$ ，如下所示：

$$y = v_i t \sin \theta - \frac{1}{2} g t^2 = 0$$

$$\rightarrow t = \frac{2v_i \sin \theta}{g}$$

$$\rightarrow v_i = \frac{gt}{2 \sin \theta}$$

$$\rightarrow \theta = \sin^{-1} \left( \frac{gt}{2v_i} \right)$$

另外，这还将简化计算  $x$  的公式：

$$x = v_i \left( \frac{2v_i \sin \theta}{g} \right) \cos \theta$$

$$\rightarrow x = \frac{2v_i^2 \sin \theta \cos \theta}{g}$$

使用三角恒等式  $\sin \theta \cos \theta = \frac{1}{2} \sin 2\theta$ ，可以将上述公式进一步简化为：

$$x = \frac{v_i^2 \sin 2\theta}{g} \quad (2.4.5)$$

另外，对于前面地面大炮向平坦地形上的地面目标开火的情况，可以使用上述公式来计算仰角为  $\theta$ 、最大初始速度为  $v_{\max}$  时，大炮的最大水平射程：

$$\text{Range} = \left| \frac{v_{\max}^2 \sin 2\theta}{g} \right| \quad (2.4.6)$$

## 2. 计算仰角

有了导弹的运动方程并解决特殊情况后，可以解决更普通的情况。首先分析发射点和目标点的海拔不同的情况，我们必须计算为将导弹从发射点发射到目标点，导弹的仰角和初始速度。由于我们已经使用  $t$  表示  $x$  和  $y$  的公式，因此可以消除其中的  $t$ ，用  $x$  和  $y$  彼此表示对方。

$$x = v_i t \cos \theta \rightarrow t = \frac{x}{v_i \cos \theta}$$

接下来，将表示  $y$  的公式中的  $t$  替换为  $x/v_i \cos \theta$ ，以消除该公式中的  $t$ 。

$$y = v_i t \sin \theta - \frac{1}{2} g t^2$$

$$\rightarrow y = v_i \left( \frac{x}{v_i \cos \theta} \right) \sin \theta - \frac{1}{2} g \left( \frac{x}{v_i \cos \theta} \right)^2$$

$$\rightarrow y = x \tan \theta - \frac{gx^2}{2v_i^2 \cos^2 \theta}$$

然后使用三角恒等式做进一步的简化。

$$\begin{aligned} y &= x \tan \theta - \frac{gx^2}{2v_i^2} \left( \frac{1}{\cos^2 \theta} \right) \\ \rightarrow y &= x \tan \theta - \frac{gx^2 (\tan^2 \theta + 1)}{2v_i^2} \\ \rightarrow \frac{gx^2}{2v_i^2} \tan^2 \theta - x \tan \theta + \frac{gx^2}{2v_i^2} + y &= 0 \end{aligned} \quad (2.4.7)$$

虽然上述公式看起来有些怪异，但它是一个二次方程式，可以使用下面的公式来求解：

$$\tan \theta = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

其中：

$$\begin{aligned} a &= \frac{gx^2}{2v_i^2} \\ b &= -x \\ c &= \frac{gx^2}{2v_i^2} + y \end{aligned}$$

将 a、b 和 c 代入到上述二次方程式的解中，可得到下面的公式：

$$\begin{aligned} \theta &= \tan^{-1} \left( \frac{x \pm \sqrt{x^2 - 4 \left( \frac{gx^2}{2v_i^2} \left( \frac{gx^2}{2v_i^2} + y \right) \right)}}{\frac{gx^2}{v_i^2}} \right) \\ \rightarrow \theta &= \tan^{-1} \left( v_i^2 \frac{x \pm \sqrt{x^2 - \frac{g^2 x^4}{v_i^4} + \frac{2gx^2 y}{v_i^2}}}{gx^2} \right) \end{aligned} \quad (2.4.8)$$

在给定初始速度  $v_i$ 、水平距离  $x$  和垂直距离  $y$  的情况下，上述二次方程式提供了  $\theta$  的解。如果  $(b^2 - 4ac)$  为正，则有两个解；如果为负，则没有解。另外，如果初始速度为 0，则弹道将垂直向下，因此  $\theta$  将无关紧要。

存在两个可能的弹道时, 扁平的弹道射向目标的速度更快, 因此大多数情况下, 这种弹道更佳。如果两个发射角都位于  $-\pi/2$  到  $\pi/2$  之间, 则在  $v_i$  相同的情况下, 更接近 0 的发射角获得的弹道更扁平。图 2.4.2 显示了两个有效的仰角, 它们都将射向给定的目标, 其中弹道 2 更早命中目标。

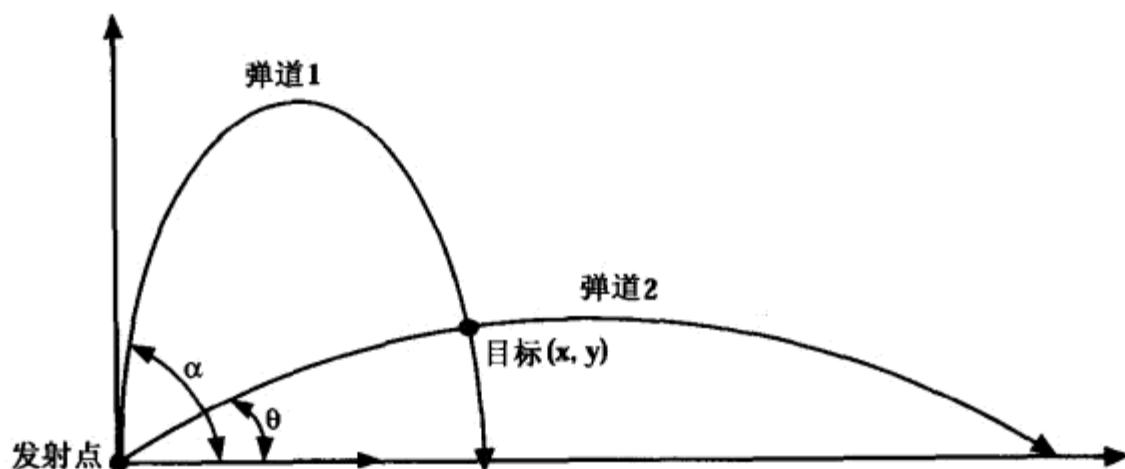


图 2.4.2 对于给定的  $v_i$ , 存在两个有效地发射角

### 3. 计算初始速度

前面介绍了如何计算  $\theta$ , 我们可以对方程 2.4.7 进行变换, 以便在发射角  $\theta$ 、水平距离  $x$  和垂直距离  $y$  已知的情况下, 计算出初始速度。

$$\frac{gx^2}{2v_i^2} \tan^2 \theta - x \tan \theta + \frac{gx^2}{2v_i^2} + y = 0$$

$$\rightarrow \frac{gx^2}{2v_i^2} \tan^2 \theta + \frac{gx^2}{2v_i^2} = x \tan \theta - y$$

$$\rightarrow \frac{gx^2}{2v_i^2} (\tan^2 \theta + 1) = x \tan \theta - y$$

将上述方程两边同时乘以  $v_i^2/(x \tan \theta - y)$ , 以分离出初始速度:

$$\frac{gx^2}{2(x \tan \theta - y)} (\tan^2 \theta + 1) = v_i^2$$

$v_i$  的解如下:

$$v_i = x \sqrt{\frac{g(\tan^2 \theta + 1)}{2(x \tan \theta - y)}} \quad (2.4.9)$$

同样, 可以使用三角恒等式  $1/\cos^2 \theta = \tan^2 \theta + 1$  对上述平方根做进一步的简化:

$$v_i = \frac{x}{\cos \theta} \sqrt{\frac{g}{2(x \tan \theta - y)}} \quad (2.4.10)$$

同样，由于解为平方根，因此有些情况下将没有解。一个这样的例子是，弹道在发射点的斜率小于目标点的斜率。一个特例是  $\theta=\pi/2$ （垂直向上发射）。

#### 4. 计算弹道的最大高度

计算弹道的最大高度很简单：最大高度的定义是：垂直速度  $v_y=0$  的点（假设  $\theta\geq 0$ ）。因此，我们只需解下述关于垂直速度的方程：

$$v_y(t) = v_i \sin \theta - gt = 0$$

$t$  的解如下：

$$t = \frac{v_i \sin \theta}{g}$$

接下来计算最大高度，用上述解替换垂直位置方程中的  $t$ ：

$$\begin{aligned} y_{\max} &= v_i t \sin \theta - \frac{1}{2} g t^2 \\ &= \frac{v_i^2 \sin^2 \theta}{g} - \frac{g}{2} \left( \frac{v_i \sin \theta}{g} \right)^2 \\ &\rightarrow y_{\max} = \frac{v_i^2 \sin^2 \theta}{2g} \end{aligned} \quad (2.4.11)$$

正如前面指出的，这要求  $\theta\geq 0$ 。如果发射角  $\theta$  小于 0（向下发射），则最大高度将为 0，因为由于地球引力，导弹的垂直速度将越来越大。从某种程度上说，这是一种特例，因为在这种情况下，导弹在最高位置时，其垂直速度比一定为 0。

#### 5. 计算飞行时间

要计算命中目标所需的时间，可以对用  $t$  表示的水平位置方程 2.4.3 进行简单的变换：

$$x(t) = v_i t \cos \theta \rightarrow t = \frac{x}{v_i \cos \theta}$$

然而，如果  $v_i=0$  或  $\cos \theta=0$ ，则用  $x$  表示时， $t$  将是未被定义的。另外，在这种情况下， $x$  总为 0，因此如果发射点和目标点的  $x$  值不同，则在这种情况下将没有解。在实现中，需要对这些边界情况进行检测，因为这种错误将导致引擎崩溃或不稳定。

要在  $v_i=0$  或  $\cos \theta=0$  的情况下计算  $t$  的值，可以使用垂直位置方程 2.4.4。

$$y(t) = v_i t \sin \theta - \frac{1}{2} g t^2$$

如果  $v_i=0$ ，则可以通过下面的公式，用  $y$  和  $g$  来表示  $t$ 。

$$y = -\frac{1}{2} g t^2 \rightarrow t = \sqrt{\frac{-2y}{g}}$$

然而，如果  $\cos \theta=0$  且  $v_i>0$ ，则可能存在一个或两个解（仅当  $\theta\geq 0$  时，才可能存在两个



解, 因为  $v_i \geq 0$ )。另外, 如果  $\cos \theta = 0$ , 则  $\sin \theta = \pm 1$ 。这将进一步简化问题, 但我们仍需要使用  $t$  来表示这一点, 如下所示:

$$\begin{aligned} y &= v_i t \sin \theta - \frac{1}{2} g t^2 \\ \rightarrow \frac{1}{2} g t^2 - v_i t \sin \theta + y &= 0 \end{aligned} \quad (2.4.12)$$

这是一个关于  $t$  的二次方程式, 如何解留给读者自己去完成。

## 6. 求解多个变量

正如本文开头指出的, 两个或更多个值 (通常是  $\theta$  和  $v_i$ , 因为发射点和目标点通常是已知的) 未知的情况很常见。对于多变量情况, 可能的解非常多, 为求解问题, 最快速的方法是消除其中的一些变量。最常见的情况是, 发射点、目标点和最大初始速度  $v_{max}$  是已知的, 需要计算  $v_i$  和  $\theta$ 。

消除变量以简化为单变量问题时, 必须以不过度限制可能解的方式进行消除, 这很重要。对于前面  $v_i$  和  $\theta$  未知的情况, 限制  $\theta$  可极大地减少解的数目。另一方面, 将  $v_i$  设置为  $v_{max}$  并改变  $\theta$  可消除大量的降落点。这种逻辑可用于这种问题的其他形式, 但由于主题的限制, 本文无法做更详细的阐述。

### 2.4.2 优化实现

使用代码实现前面的方程时, 通过进行一些优化, 可以极大地提高性能, 这是因为在大多数系统上, 三角函数的开销都非常高。

#### 1. 避免过度简化

推导数学计算公式时, 人们常常喜欢将公式简化为最简单的数学形式, 而不是最简单或最优的算法。例如, 计算初始速度时, 我们得到了方程 2.4.9 和 2.4.10, 如下所示:

$$\begin{aligned} v_i &= x \sqrt{\frac{g(\tan^2 \theta + 1)}{2(x \tan \theta - y)}} \\ &= \frac{x}{\cos \theta} \sqrt{\frac{g}{2(x \tan \theta - y)}} \end{aligned}$$

从数学的角度看, 后者更佳, 因为它对前者进行了化简; 然而, 在实现中, 预先计算  $\tan \theta$  的值并在第一个公式中使用两次的效率, 将比第二个公式计算  $\tan \theta$  和  $\cos \theta$  高。另外, 即使我们选择第二个公式 (且不化简成用  $\tan \theta$  表示), 将  $\cos \theta$  放在平方根号的外面意味着需要执行两次除法运算: 外面一次, 里面一次。要对此进行优化, 可将  $\cos \theta$  放到平方根号的除数中 (变为  $\cos^2 \theta$ ), 或者将  $x$  与  $1/\cos \theta$  相乘。

## 2. 少用三角函数以简化计算

不要使用正弦、余弦和正切函数，相反，使用预先生成的查找表或利用其他变量之间已知关系的效率将高得多。例如，要计算  $\tan \theta$ ，只须将  $v_y$  和  $v_x$  的初始值相除即可，因为它们分别是根据  $\sin \theta$  和  $\cos \theta$  定义的，而且很可能已经预先计算出来了。

另外，还有其他可优化的地方——这里旨在提醒读者，三角函数的计算开销非常高，而优化是非常重要的。

### 2.4.3 总结

---

高效的弹道生成可提高 AI 质量和引擎的性能。虽然推导涉及到大量的数学知识，但最后的公式却相对简单也容易理解。另外，理解公式的推导和化简后，将很容易将这些知识用于更复杂的情形（如移动靶、非垂直加速）和相关的问题中。



## 2.5 平行移动镜头

Carl Dougan

Carl.dougan@adrenalin.com

计 算游戏中的很多任务都要求当物体移动时形成合适的朝向。假设摄像机沿环形路径移动，您需要调整其朝向。您可能想根据摄像机行进的方向调整其朝向。当路径为环路时，摄像机的朝向应沿环路做相应的调整。您不希望摄像机朝向突然翻转（flip），而希望其方向的改变与路径的变化相称。平行移动镜头（parallel transport frame）可帮助实现这种“稳定”的定向。

这种技巧也可用于生成几何体。3D 建模中一种常见的操作是放样（lofting）——沿路径曲线延展（extrude）2D 形状，该形状生成的多个截面组合成一个 3D 几何体。如果 2D 形状为圆，则得到的 3D 模型是以路径曲线为中心的管道。前面关于摄像机的准则可用于计算形状的朝向——朝向应沿路径逐渐变化，不应发生不必要的翻转（twist）。

平行移动方法通过在沿曲线移动的过程中逐渐旋转坐标系统（镜头）来实现稳定性。关于镜头以前朝向的“记忆”使得不必要的翻转能够得以消除——在每一步中只做尽可能小的旋转，以便始终与曲线平行。不幸的是，为计算曲线末端处的镜头，必须从起点开始，沿曲线依次对镜头进行迭代，每一步都进行旋转。另两种常用的沿曲线调整镜头（framing）的方法是 Frenet Frame 和 Fixed Up[Eberly0]；对于曲线上任何一点，它们只需通过一次计算就能确定镜头的朝向。然而，这些方法也有局限性（caveat），这将在后面进行说明。

### 2.5.1 技术

一种相对简单的数值技术可用来计算平行移动的镜头[Glassner90]。对于任意初始镜头，将其沿曲线移动，并在每次迭代中，对其进行旋转，使之尽可能与曲线平行。

给定：

曲线  $C$

$t-1$  时的镜头  $F1$

$t-1$  时的切线  $T1$ （ $t-1$  时的速度或  $C$  的一次导数）

$t$  时的切线  $T2$

则  $t$  时的镜头  $F2$  可通过下述方式计算得到：

$F2$  的位置为  $t$  时  $C$  的值

$F2$  的朝向可通过将  $F1$  绕轴  $A$  旋转角度  $\alpha$  得到, 其中:

$$A = T1 \times T2$$

$$\alpha = \text{ArcCos}((T1 \cdot T2) / (|T1||T2|))$$

如果切线  $T1$  和  $T2$  平行 (即  $T1 \times T2$  为零), 则无需旋转 (如图 2.5.1 所示)。

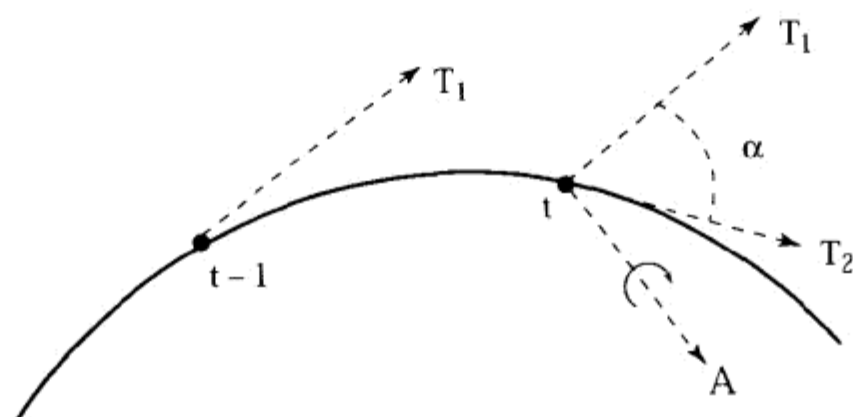


图 2.5.1 将  $t-1$  时的镜头绕轴  $A$  旋转  $\alpha$ , 从而得到  $t$  时的镜头

最初的镜头可以是任意的。可以使用 Fixed Up 或 Frenet Frame 方法来计算轴与切线重叠时的初始镜头。

有些情况下, 这样做可能是可取的, 即使用平行移动方法根据曲线上的粗略取样生成一些镜头, 然后采用四元数插值在样本镜头间实现平滑的旋转。使用四元数 (quaternion) 总是可取的, 因为存在一种高效的根据旋转轴和旋转角度生成四元数的方法[Eberly01]。您可以根据前面所示的角度和轴生成一个旋转四元数, 然后将其与前一个镜头的四元数相乘来进行旋转。

## 1. 移动的物体

对于移动的物体, 可以在物体每次移动时 (大概每帧移动一次), 通过执行一次平行移动旋转来调整其朝向。这需要三项信息: 物体在当前位置和前一个位置处的速度以及在前一个位置的朝向。速度为前图中的切线  $T1$  和  $T2$ 。

对于有些任务, 平行移动镜头可能过于“稳定”。例如, 飞机在水平面上做 S 形飞行时, 将不会倾斜。为模拟逼真的飞行, 可能需要采用其他的解决方案, 如模拟运动的物理现象。Crag Reynolds 描述了一种相对简单而快速的调整鸟群飞行朝向的技巧, 其中包含倾斜飞行 [Reynolds99]。Reynolds 的技巧与平行移动方法的相似之处在于, 它们都依赖于对前一个镜头的“记忆”。

## 2. 比较

接下来将平行移动方法同沿曲线调整镜头的 Frenet Frame 方法和 Fixed Up 方法进行比较。

### (1) Frenet Frame

Frenet Frame 使用三条正交轴:

- 曲线的切线;

- 切线与二阶导数的叉积；
- 根据前两个向量的叉积得到的另一个向量。

用于前面讨论的情况时，Frenet Frame 存在问题，因为当二阶导数为 0 时，将无法进行镜头计算。在曲线的拐点和直线部分，二阶导数将为 0 [Hanson95]。显然，就我们的目标而言，不能计算直线部分的镜头是个大问题。另外，镜头可能由于二阶导数的变化而旋转。例如，对于 S 形曲线，二阶导数指向内部，因此上半段和下半段的指向相反。因此，对于 S 形曲线，Frenet 镜头将翻转。图 2.5.2 对此进行了说明，在二阶导数的指向发生变化的地方，几何体不是连续的，而是被隔断。如果这是一群飞行的鸟，则它们在这一点将翻个筋斗，从向上飞行转为向下飞行。

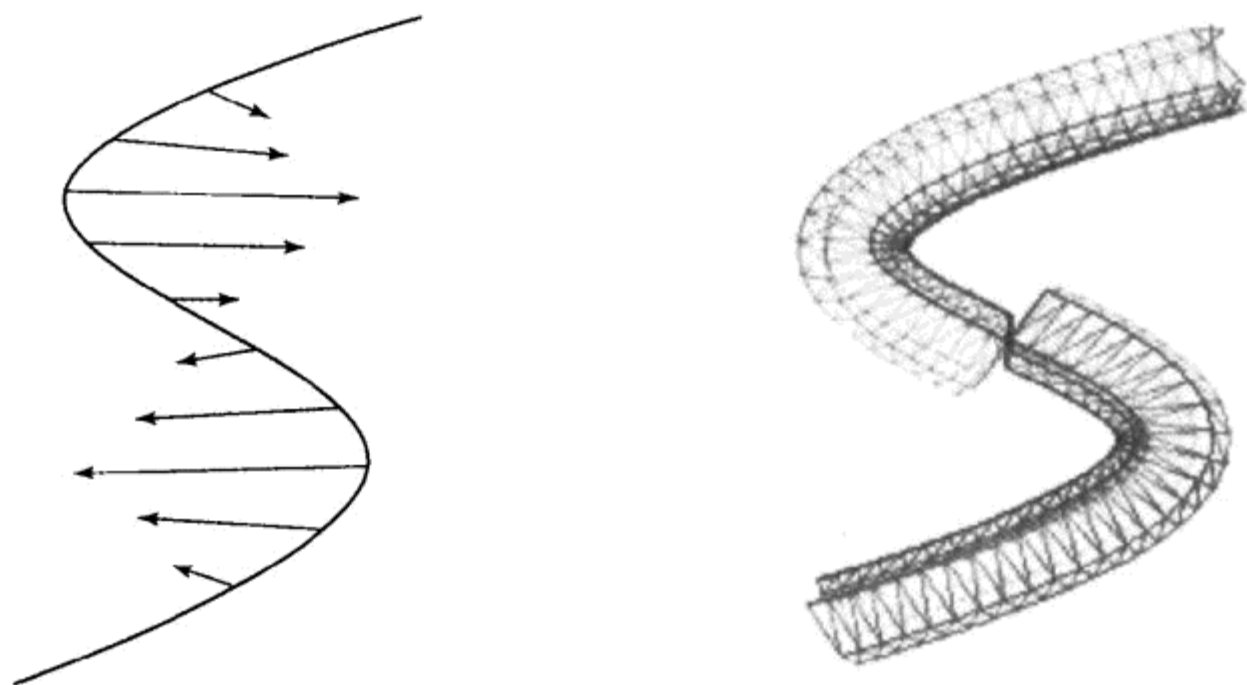


图 2.5.2 S 形曲线的二阶导数及 Frenet Frame 方法生成的曲线管

## (2) Fixed Up 方法

Fixed Up 方法使用切线  $T$  和任意向量  $V$  (Fixed Up 向量) 来生成镜头的三条轴——方向  $D$ 、向上的向量  $U$  和向右的向量  $R$  [Eberly01]:

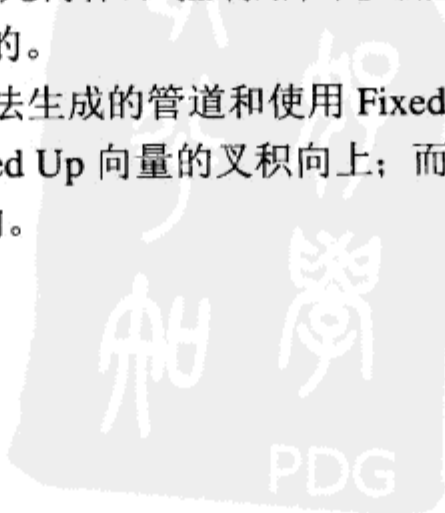
$$D = T / |T|$$

$$R = D \times V / |D \times V|$$

$$U = R \times D$$

当切线同任意选择的向量平行或几乎平行时，Fixed Up 方法将出现问题。当  $T$  和  $V$  平行时， $D$  和  $V$  的叉积为 0，镜头无法创建。即使它们几乎平行， $T$  的微小变化都将使生成的向量相对于切线发生急剧变化，导致镜头翻转。如果您能够控制路径，这将不是什么问题。对于有些任务（如构建高速公路对应的几何体），控制路径是可能的；但对于其他任务（如构建过山车对应的几何体），这是不可能的。

图 2.5.3 显示了使用平行移动方法生成的管道和使用 Fixed Up 方法生成的管道。在曲线的上半部分和下半部分，切线与 Fixed Up 向量的叉积向上；而在中间部分，则向下。这种急剧的翻转导致生成的几何体是断裂的。



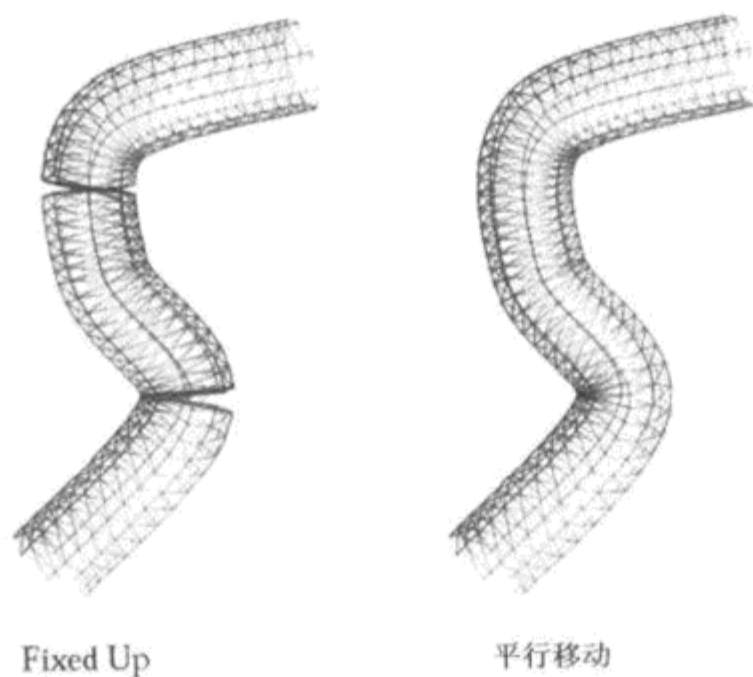


图 2.5.3 Fixed Up 方法和平行移动方法之比较

## 2.5.2 结论

对于不受约束的路径，如飞行的导弹或环形铁轨，可使用平行移动方法来避免轨道断裂和导弹翻筋斗。

## 2.5.3 参考文献

- [Glassner90] Bloomenthal, Jules, "Calculation of Reference Frames Along a Space Curve," *Graphics Gems*, Academic Press, 1990: pp. 567–571.
- [Eberly01] Eberly, David H., *3D Game Engine Design*, Academic Press, 2001.
- [Hanson95] Hanson, Andrew J., and Ma, Hui, *Parallel Transport Approach to Curve Framing*, Department of Computer Science, Indiana University, 1995.
- [Reynolds99] Reynolds, Craig, "Steering Behaviors for Autonomous Characters," available online at [www.red3d.com/cwr/steer/gdc99/index.html](http://www.red3d.com/cwr/steer/gdc99/index.html).



## 2.6 平滑的基于四元数的 $C^2$ 飞行路径

Alex Vlachos, ATI Research; John Isidoro  
avlachos@ati.com; jisodoro@cs.bu.edu

本文介绍一种对摄像机的位置和朝向进行平滑插值，以生成  $C^2$  连续的飞行路径 (flythrough path) 的方法。我们吸取了多种著名方法的精华，并提供了实现本文中介绍的方法的 C++ 类。

### 2.6.1 导论

通过对样本点应用自然三次样条函数 (natural cubic spline)，可以很容易地对飞行路径的位置进行插值；然而朝向则要复杂些。我们将描述一种将  $S^3$  空间中的四元数 (单位超球面上的点) 转换到  $R^4$  空间 (4D 空间中的点) 的方法 [Johnstone99]。四元数位于  $R^4$  空间中后，便可以将任何 4D 函数用于转换后的数据。然后，可以将通过插值得到的点转换到  $S^3$  空间中，并作为四元数使用。另外，还将描述一种名为选择性求负 (selective negation) 的技术，它对四元数进行预处理，得到样品点朝向之间的最短旋转路径。

镜头剪接 (camera cut，将摄像机移到另一个位置) 是通过在剪接的前后加上幻影点 (phantom point) 来实现的，就像插入非闭合样条线那样。



为插入样条线，以便在剪接点附近获得平滑的结果，这些点是必不可少的。附带光盘中的代码将剪接点作为飞行路径的一部分，这样代码更简单。在 C++ 类的内部，剪接段被视为独立的样条线，没有为每段创建样条线的开销。

### 2.6.2 位置插值

下面讨论位置插值。

#### 1. 样本点

常见的指定样本点的方式有两种。第一种方式是，控制点之间的每条线段表示的时间相同 (如 1 秒)；第二种方式是，控制点只用于定义摄像机路径的形状，而摄像机以恒定的速度沿路径前进。本文提供的代码假设控制点之间的时间是恒定的，当然可以很容易地对这些代码进行修改，用于





示旋转所需的信息，不多也不少。

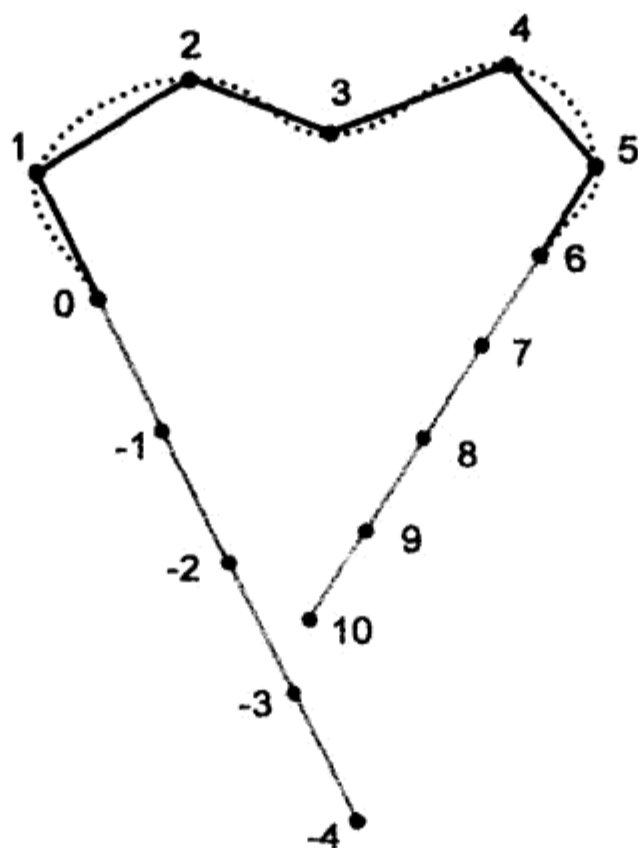


图 2.6.2 为非闭合样条线创建幻影点

然而，对于朝向而言，使用四元数存在多义性。朝向可被视为在已有朝向的基础上进行旋转。使用四元数时，存在两个得到相同朝向的旋转。假设绕轴  $w$  逆时针旋转  $\theta$  可得到所需的朝向，则绕轴  $-w$  旋转  $360^\circ - \theta$  也将得到这样的朝向。转换为四元数后，第二个四元数是第一个的负（negation）。

#### 2.6.4 旋转方向 and 选择性求负

进行四元数插值时，有一个微妙之处需要注意。表示朝向时，一个四元数或其负便足够了。然而，进行朝向插值时（如执行旋转），正四元数和负四元数将导致完全不同的旋转，进而导致不同的摄像机路径。如果希望在每对朝向之间执行尽可能小的旋转，则可以通过对四元数进行预处理来实现。

两个四元数的点积是它们之间的旋转角度的一半的余弦。如果这个值为负，则说明这两个四元数之间的旋转角度大于  $180^\circ$ 。在这种情况下，对其中一个四元数求负，将使旋转角度小于  $180^\circ$ 。对插值而言，这使得在关键镜头的朝向之间总是进行最短的旋转。我们称这种过程为选择性求负。

对朝向四元数进行选择性求负技术可作为摄像机飞行路径的预处理步骤。为执行预处理步骤，沿飞行路径从头走到尾，对于路径上的每个四元数  $q_i$ ，如果它与前一个四元数的点积为负（即  $(q_i \cdot q_{i-1}) < 0$ ），则对其进行求负。使用选择性求负作为预处理步骤，使得样条线插值的效率高得多，因为无需对每个样本进行选择性求负。

要对闭合的样条线路径进行预处理，必须在选择性求负之前复制样条线的前 4 个点，并

将它们加入到路径的末尾。注意，复制得到的点的符号可能与原来的点不同。对于非闭合样条线，则需要创建一些幻影四元数（对应于幻影控制点），并将其加入到样条线两端。这类似于需要在离路径两端最近的两个四元数之间进行线性插值。然而，以线性方式插入四元数还不够。相反，我们将使用球状线性（spherical linear）算法。给定四元数  $q_0$  和  $q_1$ ，我们需要生成 4 个幻影四元数（ $q_{-1}$ 、 $q_{-2}$  等），并将其插入到非闭合样条线的开头。我们使用球状线性插值（slerp）函数来在  $q_0$  和  $q_1$  之间进行插值，其中使用的 slerp 值为 2.0。这样，幻影点的旋转将以线性的方式变化。

针对插值对所有的朝向四元数进行预处理后，执行平滑的基于样条线的四元数插值技术将非常简单。

### 2.6.5 四元数样条线插值

正如我们在位置插值中指出的，样条线可用于实现比线性插值更平滑的插值。然而，四元数样条线插值并非那么简单，可供使用的技术有多种。一种技术是，根据未经处理的四元数进行插值，然后对得到的四元数重新进行规格化。然而，这种技术得到的路径并非平滑的，且会给角速度带来奇怪的变化。另一种想法是，使用基于四元数对数的技术。一个这样的例子是 SQUAD（球状四边形插值，spherical quadrangle interpolation）。使用这些技术时，性能将受到限制，因为它们需要使用超越函数（正弦、余弦、对数、幂等）。其他技术涉及到混合单位四元数超球面上的两个球[Kim95]或某种迭代数值技术[Neilson92]。虽然很多这样的技术都能得到相当好的结果，但大多数都不具备  $C^2$  连续性或者计算量过大无法采用，尤其是使用大量的飞行路径时（例如，游戏角色或导弹的飞行路径）。

然而，有一种四元数样条线插值技术能够得到非常好的结果，并能满足导数连续性要求。这种技术在单位四元数 4-球（ $S^3$ ）和另一个四维空间（ $R^4$ ）之间使用一个可逆的有理映射 [Johnstone99]M。在下面的方程中， $a$ 、 $b$  和  $c$  是四元数的向量部分的分量，而  $s$  为标量部分。

从  $S^3$  到  $R^4$  的变换  $M^{-1}$  为：

$$x = a / \sqrt{2(1-s)}$$

$$y = b / \sqrt{2(1-s)}$$

$$z = c / \sqrt{2(1-s)}$$

$$w = (1-s) / \sqrt{2(1-s)}$$

从  $R^4$  到  $S^3$  的变换  $M$  为：

$$s = (x^2 + y^2 + z^2 - w^2) / (x^2 + y^2 + z^2 + w^2)$$

$$a = 2xw / (x^2 + y^2 + z^2 + w^2)$$

$$b = 2yw / (x^2 + y^2 + z^2 + w^2)$$

$$c = 2zw / (x^2 + y^2 + z^2 + w^2)$$

将这种技术用于四元数样条线插值很简单。首先，必须对控制四元数进行选择性求负，确保在控制点之间进行最短的旋转。然后，对所有的控制四元数应用  $M^{-1}$ ，以获得它们在  $R^4$  中对应的值。这可以作为一个预处理步骤来完成，也可以在构建飞行路径时或程序的装载阶段完成。这样，在计算飞行路径时，无需计算平方根。

接下来，使用您选择的样条线对得到的四元素向量进行插值。由于这是一个连续的有理

映射, 因此通过插值得到的  $S^3$  四元数路径的连续性与用于在  $R^4$  空间进行插值的样条线相同。

在我们的应用中, 我们使用自然三次样条线[Hearn94][Press97]在  $R^4$  空间中进行插值。

这样朝向插值的连续性为  $C^2$ 。从质量上说, 这样做的效果是, 摄像机路径的交角速度不会发生急剧变化。

找到  $R^4$  空间中所需的样条线上的点后, 便可以使用  $M$  将其转换为四元数。

然而, 使用这种技术时, 还有一个数学上的微妙问题需要解决。

### 2.6.6 有理映射中的奇异点

如果飞行路径中包含等于或接近于  $(1,0,0,0)$  的朝向四元数, 则通过  $R^4$  空间进行插值时, 将导致数值不稳定性, 因为  $M^{-1}(1,0,0,0) = (\infty, \infty, \infty, \infty)$ 。对于这种奇异性, 处理的方式有几种。一种方式是忽略它, 令人惊讶的是, 在很多情况下, 这都是可行的。例如, 如果  $z$  轴垂直向上, 而您知道摄像机不会垂直向上, 导致摄像机的  $up$  向量指向  $y$  轴, 因此朝向四元数永远不会等于  $(1,0,0,0)$ , 这样问题就解决了。

如果情况并非如此, 则另一种办法是, 找出一个四元数  $q_f$ , 它与任何朝向四元数之间的角度都大于  $30^\circ$ , 然后使用  $q_f$  将所有的四元数都旋转为远离奇异点的安全朝向[Johnstone99]。这里的基本思想是, 对旋转后的飞行路径进行插值, 然后将插值得到的朝向旋转到原来的坐标系中。这里需要做的只是在执行选择性求负步骤后, 将所有的朝向四元数乘以  $q_f$ 。接下来, 将四元数从  $S^3$  空间变换到  $R^4$  空间、应用自然三次样条线, 并将得到的  $R^4$  值变换到  $R^3$  空间。然后, 执行多出的一步, 在使用每个四元数之前, 将其旋转  $q_f^{-1}$ 。

一个查找  $q_f$  的简易方式是, 随机的生成单位四元数, 直到出现一个与任何一个被选择性求负的朝向四元数之间的角度都大于  $30^\circ$  的单位四元数为止。

### 2.6.7 镜头剪接

镜头剪接的定义是, 将摄像机从场景中的一个点移向另一点。不能在样条线中间插入剪接, 也不能跨过样条线中的一段。相反, 需要以剪接点为界将样条线分为两条, 并将它们作为独立的非闭合样条线单独进行处理。我们提供的代码以这种方式来处理镜头剪接, 即无需显式地创建一条独立的路径 (如图 2.6.3 所示)。

### 2.6.8 代码



本文附带的代码是一个 C++ 类, 它实现了正文中介绍的大多数技术。这个类包含用于手工创建和编辑控制点、读写路径文件、处理剪接点、对给定时间的样条线值进行取样以及建立顶点和索引缓冲区以便绘制样条线的成员函数。这个类假设每个控制点之间相隔的时间相同, 而不是假设摄像机的前进速度恒定。另外, 这里的代码没有解决奇异点的问题, 因为在我们的项目中从来没有发生过这种奇异现象。更详细的信息, 请参阅源代码文件。

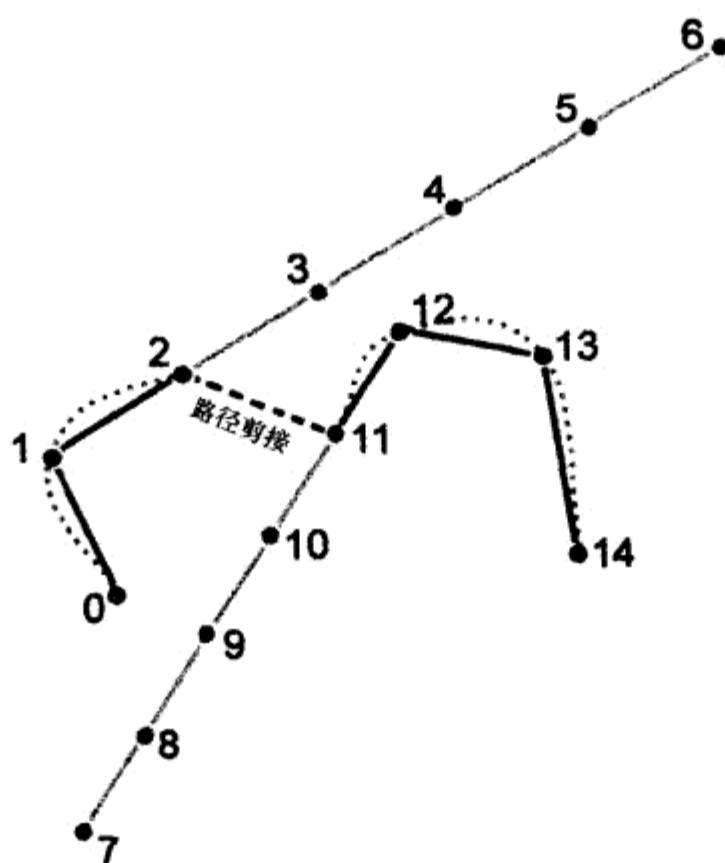


图 2.6.3 为路径剪接生成幻影点

## 2.6.9 参考文献

[Press97] Press, William H., et al, *Numerical Recipes in C*, Cambridge University Press, 1997.

[Hearn94] Hearn, Donald, Baker, M. Pauline, *Computer Graphics Second Edition*, Prentice Hall, Inc. 1994.

[Johnstone99] Johnstone, J. K., Williams, J. P., "A Rational Quaternion Spline of Arbitrary Continuity," Tech Report: [www.cis.uab.edu/info/faculty/jj/cos.html](http://www.cis.uab.edu/info/faculty/jj/cos.html), 1999.

[GPG1] Edited by Mark DeLoura, *Game Programming Gems*, Charles River Media, 2000.

[Shoemake91] Shoemake, K., *Quaternion Calculus for Animation*, Math for SIGGRAPH (ACM SIGGRAPH '91 Course Notes #2), 1991.

[Neilson92] Neilson, G., and Heiland, R., "Animating Rotations Using Quaternions and Splines on a 4D Sphere," English Edition, *Programming and Computer Software*, Plenum Pub., New York. 1992.

[Kim95] Kim, M.S. and Nam, K.W., *Interpolating Solid Orientations with Circular Blending Quaternion Curves*, Computer-Aided Design, Vol. 27, No. 5, pp. 385–398, 1995.

## 2.7 递归逐维分组：一种快速的碰撞检测算法

Steve Rabin, 任天堂 (美国) 公司  
steve\_rabin@hotmail.com

**基**本的碰撞检测 (collision detection) 是一种成本高昂的操作。这种最简单的方法将每个物体视为一个球，需要对距离进行  $n^2$  次比较。程序清单 2.7.1 是这种简单、蛮干的比较算法的代码。然而，对于典型的情况，一种名为递归逐维分组 (Recursive Dimensional Clustering, RDC) 的新技术可以将复杂度从  $O(n^2)$  降低到  $O(n\log^2 n)$ 。

程序清单 2.7.1 蛮干的比较算法

```
// Slightly optimized  $O((n^2-n)/2)$  time complexity
for( i=0; i<num_objects; i++ ) {
    for( j=i+1; j<num_objects; j++ ) {
        if( Distance(i, j) < Radius(i) + Radius(j) ) {
            //in collision
        }
    }
}
```

表 2.7.1 列出了 RDC 方法和蛮干方法的差别。其中的时间是使用只包含一次碰撞的数据在 Pentium II 400 MHz 的机器上得到的。注意，程序清单 2.7.1 对蛮干的算法进行了最大限度的优化，将其复杂度从  $O(n^2)$  降低到了  $O((n^2-n)/2)$ ，同时比较时使用距离的平方，从而省却了计算平方根的成本。

表 2.7.1 RDC 方法和蛮干方法之间的差别

物体数目	RDC 方法	蛮干方法
10	< 1 ms	< 1 ms
50	2 ms	6 ms
100	4 ms	18 ms
200	7 ms	67 ms
300	11 ms	150 ms
400	15 ms	263 ms
500	19 ms	406 ms
1000	38 ms	1596 ms (1.6s)
2000	81 ms	6335 ms (6.3 s)
5000	222 ms	39380 ms (39.4 s)
10000	478 ms	157621 ms (2 min 37.6 s)



表 2.7.1 有趣的一点是, 对于 RDC 方法, 时间与物体数目呈线性关系; 而对于蛮干的方法, 时间与物体数目呈指数关系。虽然改进的效果非常显著, 但今当处理大量通常不会发生碰撞的物体时, RDC 才有用武之地。

### 2.7.1 其他应用

RDC 的优点之一是, 它能够找出碰撞组, 而不是两个碰撞的物体。组指的是一系列彼此接触或彼此的距离在某个半径范围内的物体。该定义还包含这样的含义, 即如果 A 与 B 相接触, 而 B 又与 C 相接触, 则 A、B 和 C 属于同一组。这意味着组成员并不一定要彼此都接触。因此, RDC 可用来检测多个物体同时相撞的情况, 甚至战场上多对敌军相遇的情况。

RDC 一个有趣的应用是, 识别元球 (metaball cluster) 以提高几何体的生成速度 (Marching Cubes 算法)。在这种应用中, RDC 找出元球簇, 然后计算这些元球簇的最小同轴 (axis-aligned) 外接长方体。最小外接长方体 (bounding box) 对 Marching Cubes 算法来说至关重要, 这是因为计算速度与外界长方体的体积成正比, 因此这种算法的时间复杂度为  $O(n^3)$ 。第 3 章中的第一个彩页显示了一个例子。

#### 1. RDC 算法

要理解 RDC 的原理, 首先需要知道它如何沿一维识别碰撞物体组。图 2.7.1 是一个一维的例子, 其中有 3 个物体。

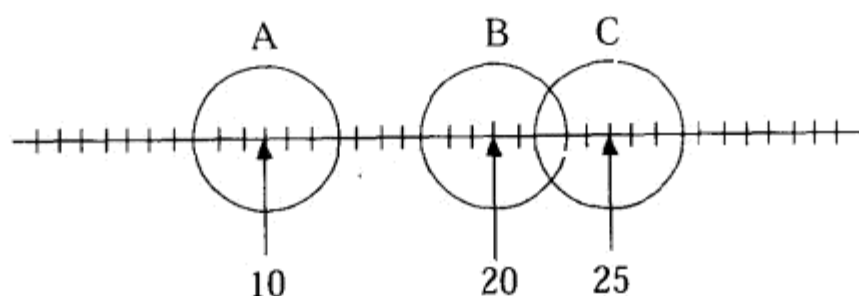


图 2.7.1 一维空间中的三个物体

正如您看到的, 物体 B 和 C 重叠在一起, 而物体 A 单独放置。因此分组 (clustering) 算法将认为有两组: 一组包含 A, 另一组包含 B 和 C。虽然通过视觉确定这一点很容易, 但我们需要一种能够得到相同结果的系统算法。

这个算法的基本思想是, 对每个物体的边界进行标记, 然后使用这些数据找出边界重叠的物体组。这是通过遍历所有的实体, 并将其在给定维中的起始边界和结束边界存储在一个链表中实现的。例如, 物体 A 的左边界为 7, 右边界为 13。为使这种标记方式适用于任何维, 可以将这些边界看作是方括号, 并将其标记为起始 (open) 和结束 (close), 而不是左边界和右边界。图 2.7.2 是根据 3 个物体创建的链表。

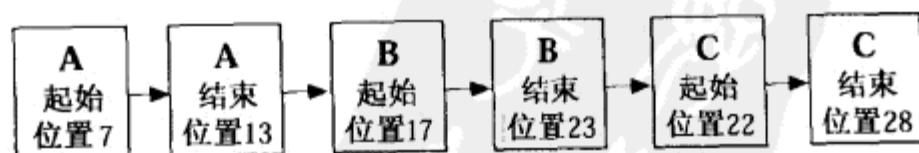


图 2.7.2 边界链表



创建边界链表后，接下来需要使用您喜欢的排序算法（如快速排序算法）按位置从低到高的顺序对链表进行排序。排序后的链表如图 2.7.3 所示。在这个例子中，只有两个元素被互换。

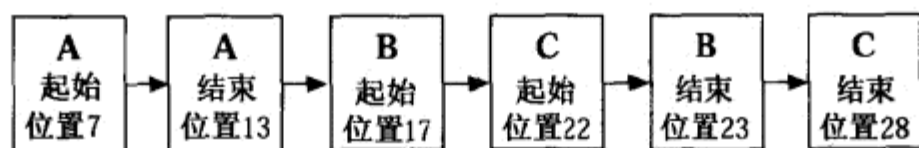


图 2.7.3 排序后的链表

有了经过排序的边界链表后，就可以找出物体组。查找物体组的算法与查找匹配方括号的分析算法非常类似。程序清单 2.7.2 中的伪代码遍历排序后的链表，并找出其中的物体组。

#### 程序清单 2.7.2 沿一维找出物体组的算法

```

int count = 0;
Clear( currentGroup );
for( element in list )
{
    if( element is an "open bracket" ) {
        count++;
        Add entity to currentGroup;
    }
    else { //element is a "closed bracket"
        count--;
        if( count == 0 ) { //entire group found
            Store( currentGroup );
            Clear( currentGroup );
        }
    }
}
assert( count == 0 );

```

至此，您可能注意到了，该算法将共享边界的物体视为一组。例如，如果物体 A 的结束边界的位置为 10，而物体 B 的开始边界的位置也为 10，则任何简单的排序算法都无法区分它们。结果是，程序清单 2.7.2 将不会以一致的方式对这种情况进行分组。然而，对于这种问题，解决方案有很多。最简单的办法是，使用浮点值，从而确保物体的边界位置不为整数。另一种解决方案是，将每个物体的半径都加大一些，这样相同的边界将彼此分离。也可以在排序函数中解决这种问题，然而，这将带来额外的开销，进而增加运行时间。

## 2. 多维空间中的 RDC

显然，在一维空间中，这种算法是管用的。然而，仅当用于二维或三维空间时，这种算法才能发挥作用。对于多维空间，窍门在于首先沿一条轴查找物体组，然后沿下一条轴进一步对新的组进行分组。

图 2.7.4 是一个二维空间的例子，其中包含 4 个物体。通常，根据观察，很容易对它们进行分组，但这个例子将演示这种算法如何对多维空间中的物体进行分组。

对于图 2.7.4 中的物体，沿  $x$  轴的排序后的链表如图 2.7.5 所示。

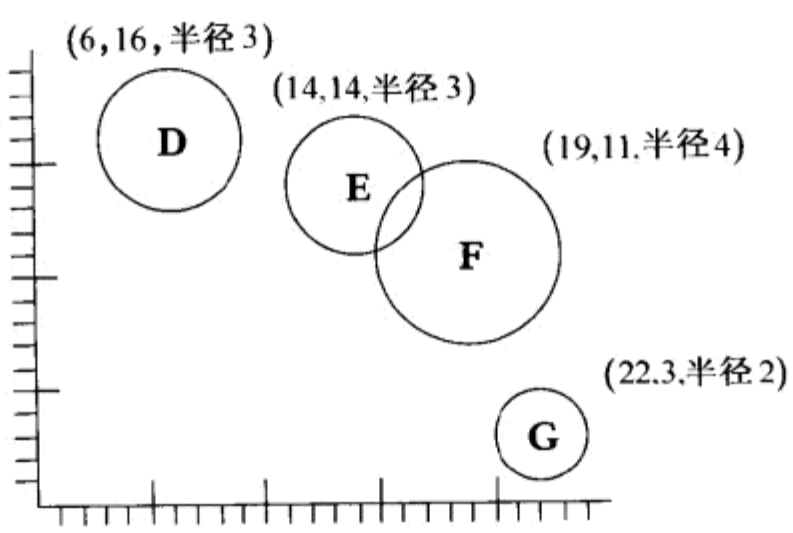


图 2.7.4 包含 4 个物体的二维范例

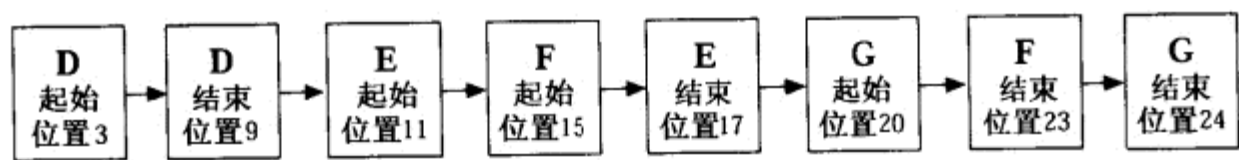


图 2.7.5 图 2.7.4 对应的排序后的边界链表

沿  $x$  轴分组得到的结果如下：

$Group_0 = \{ D \}$ ,  $Group_1 = \{ E, F, G \}$

至此，我们知道至少有两组。第一组包含物体 D，第二组包含物体 E、F 和 G。沿  $x$  进行分组后，对每个新的组都沿  $y$  轴进一步进行分组。由于物体 D 所在的组只有一个成员，因为无需在进行分组。然而，现在必须沿  $y$  轴对物体 E、F 和 G 进行分析。图 2.7.6 是  $Group_1$  沿  $y$  轴的排序后的链表。

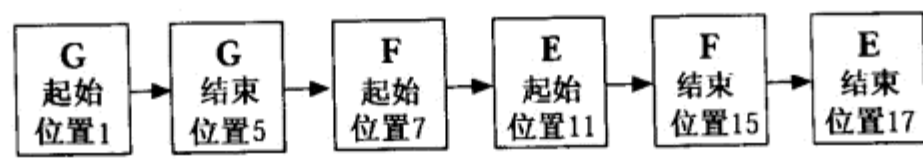


图 2.7.6  $Group_1$  沿  $y$  轴的排序后的边界链表

沿  $y$  轴进行分组的结果如下：

$Group_{1a} = \{ G \}$ ,  $Group_{1b} = \{ E, F \}$

至此，我们沿每一维都进行了分析，最终结果如下：

$Group_0 = \{ D \}$ ,  $Group_{1a} = \{ G \}$ ,  $Group_{1b} = \{ E, F \}$

图 2.7.7A 和 2.7.7B 说明了每个物体在各维中的边界以及最后的分组结果。要将这种算法扩展到三维，只需再沿  $z$  轴对各组进行分析即可。现对 RDC 算法的步骤总结如下：

- (1) 从  $x$  轴开始；
- (2) 创建这一维的物体边界链表；
- (3) 按边界位置从低到高的顺序对链表进行排序；
- (4) 使用程序清单 2.7.2 所示的开始/结束“方括号匹配”算法找出物体组；

(5) 对于找到的每个分组，沿下一个坐标轴重复第 (2) ~ (5) 步。

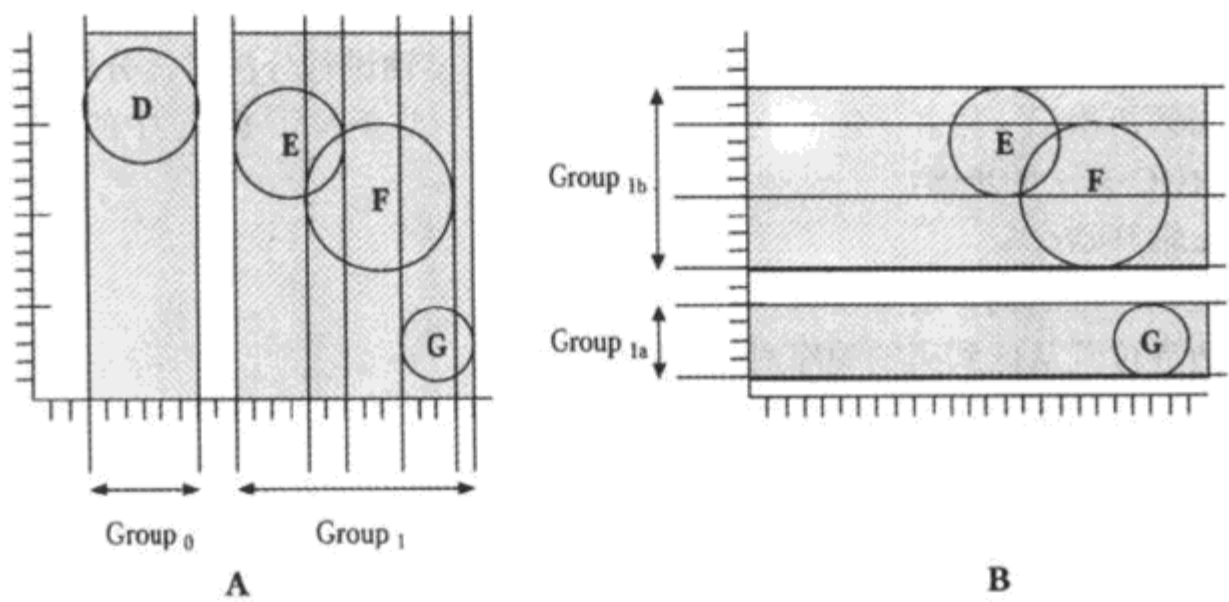


图 2.7.7 A) 沿  $x$  轴找出的分组; B) 沿  $y$  轴对  $\text{Group}_1$  进行分组

2.7.2 该算法的缺陷

不幸的是，前面描述的算法有个致命的缺陷。沿  $x$  轴对物体进行分组时，该算法可能被假象迷惑将物体分为一组，而沿其他轴无法正确地将该组物体进行划分。图 2.7.8 说明了这样一种情况。

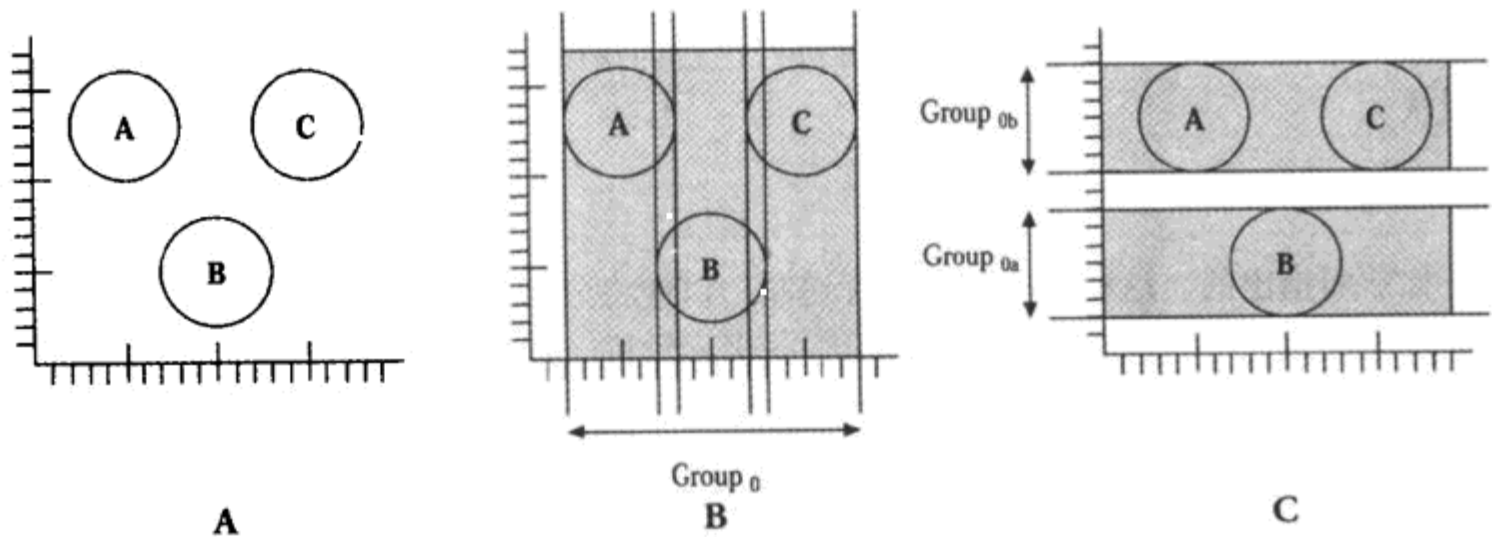


图 2.7.8 A) 算法出现缺陷的情况; B) 沿  $x$  轴分组; C) 沿  $y$  轴进一步划分  $\text{Group}_0$

图 2.7.8B 表明，沿  $x$  轴进行分组时，发现所有 3 个物体重叠在一起，因此它们被划分为一组。然后，沿  $y$  轴对该组进行分析，如图 2.7.8C 所示。其结果是，这些物体被划分为两组，而物体 A 和 C 没有被划分到不同的组。

正确的答案应该是，每个物体都属于单独的一组。然而，该算法的结果不完全正确。对于这种问题，解决的办法是再次沿  $x$  轴对沿  $y$  轴找出的各个组进行分析。这样，便可以找到正确的分组。

可将该解决方案进行推广，以确保该算法总是正确。该算法应遵循下述规则：

沿某个坐标轴对分组进行进一步划分后，必须再次沿其他所有坐标轴对得到的子组进行

分析。

就二维而言，如果沿  $y$  轴进行分析时，分组被进一步划分为子组，则必须再次沿  $x$  轴对这些子组进行分析，反之亦然。在三维空间中，如果沿  $y$  轴进行分析时，分组被进一步划分为子组，则必须再次沿  $x$  轴和  $z$  轴对这些子组进行分析。这解释了 RDC 中的“递归”。下面是修订后的 RDC 算法的步骤：

- (1) 从  $x$  轴开始；
- (2) 创建这一维的物体边界链表；
- (3) 按边界位置从低到高的顺序对链表进行排序；
- (4) 使用程序清单 2.7.2 所示的开始/结束“方括号匹配”算法找出物体组；
- (5) 对于找到的每个分组，沿其他坐标轴重复第 (2) ~ (5) 步，知道给分组不被进一步划分，且沿所有坐标轴对该分组进行了分析。

### 2.7.3 查找碰撞物体对

正如前面介绍的，RDC 只识别相接触的一组物体。物体组包含的成员可能两两彼此直接相接触，也可能不是这样。虽然这有很多用途（同时碰撞、分组信息），但通常的碰撞检测要求找出相碰撞的两个物体。

要找出碰撞的物体对，必须将 RDC 算法得到的每一组传递给蛮干比较算法。此时，组中的物体数目很可能很少，因此这种复杂度为  $O((n^2-n)/2)$  的算法的速度将足够快。

一种可能更快地找出碰撞物体对的方式是，一旦碰撞组中的物体数目小于某个值（如 10），便采用蛮干算法。此时，对所有物体进行比较的速度，将比继续通过递归调用进一步划分分组快。

#### RDC 实现

正如前面指出的，这种算法不断进行递归，以便将分组划分为最小的簇。这里比较棘手的地方是：设置一个递归函数，它选择沿哪些坐标轴进一步对分组进行划分，以及在什么情况下结束递归。在三维空间中，函数必须至少沿所有 3 个坐标轴进行划分。然而，如果沿某个坐标轴进行分析时，导致了进一步的分组，则必须递归地沿其他两个坐标轴对这些子组进行分析。

为完成这种任务，最简单的办法是，让递归函数接受 3 个指出还需要分析哪些维的参数。当该函数首次被调用时，所有三维都出现在参数列表。对每维进行分析时，如果没有进一步分组，则维参数列表将缩短为 0，递归结束；如果进一步分组了，则使用两个参数（另外两维）递归调用函数。完整的函数如程序清单 2.7.3 所示。

#### 程序清单 2.7.3 RDC 算法（伪代码）

```
void RDC( Pairs& pairs, Group& group,
          Axis axis1, Axis axis2, Axis axis3 )
{
    // "pairs" holds all final pairs that are in collision
    // "group" is the current group of objects to analyze
```

```
    // "axis1" is the axis to analyze within this function
    // "axis2", "axis3" will be analyzed in recursive calls

    if( Size( group ) < 10 || axis1 == INVALID_AXIS )
    { //end recursion and test for collisions
        BruteForceComparison( pairs, subGroup );
    }
    else {
        //for this group, find the boundaries and sort them
        OpenCloseBoundaryList boundaries;
        FindOpenCloseBoundaries( axis1, group, boundaries );
        SortOpenCloseBoundaries( boundaries ); //O(nlogn)

        Group subGroup;
        unsigned int count = 0;
        Axis newAxis1 = axis2;
        Axis newAxis2 = axis3;
        Axis newAxis3 = INVALID_AXIS;
        bool groupSubdivided = false;

        //subdivide the group if possible and call recursively
        for( every curBoundary in boundaries list )
        {
            if( curBoundary is "open bracket" )
            { //this entity lies within a cluster group
                count++;
                AddToGroup( subGroup, curBoundary->entity );
            }
            else
            {
                count--;
                if( count == 0 )
                { //found the end of a cluster group - take subgroup
                    //and call recursively on remaining axis'

                    if( curBoundary != GetLastBoundary( boundaries ) )
                    { //this group is being subdivided - remember
                        groupSubdivided = true;
                    }

                    if( groupSubdivided )
                    { //reconsider all other axis'
                        if( axis1 == X_AXIS ) {
                            newAxis1 = Y_AXIS;
                            newAxis2 = Z_AXIS;
                        }
                        else if( axis1 == Y_AXIS ) {
                            newAxis1 = X_AXIS;
                            newAxis2 = Z_AXIS;
                        }
                    }
                }
            }
        }
    }
}
```

```

        else if( axis1 == Z_AXIS ) {
            newAxis1 = X_AXIS;
            newAxis2 = Y_AXIS;
        }
    }
}

//recursive call
RecursiveClustering( pairs, subGroup,
                    newAxis1, newAxis2, INVALID_AXIS );
Clear( subGroup ); //clear the subGroup for the next group
}
}
}
}
}

```

上述 RDC 函数还具备查找碰撞物体对的功能。递归结束并找到最小的分组后，分组将被传递给蛮干比较函数。如果分组中的物体少于 10 个，递归也将结束，并立刻使用蛮干函数对成员进行比较。

#### 2.7.4 时间复杂度

乍一看，该算法好像非常耗时。直觉告诉您，该算法的复杂度为  $O(n^2)$ ——这很接近。然而，由于三维空间的物理限制，最糟的情况出现的可能性非常小。相反，只要大部分物体没有相碰，则该算法的复杂度为  $O(n \log_2 n)$ 。

在两种极端情况下，RDC 的性能非常糟。一种情况下，物体的分布情况导致递归层次非常深，这是最糟的情况，复杂度为  $O(n^2 \log_2 n)$ 。另一种情况是，所有物体碰撞到一起。在这种情况下，RDC 几乎没有做任何工作，必须使用蛮干算法，复杂度为  $O(n^2)$ 。

在最糟的情况下，递归层次非常深。这发生在物体分布最不对称时，即其中一个物体属于一组，其他  $n-1$  物体属于另一组。这样，大的一组将被递归地传递给函数。如果每个递归层的情况都如此，则需要进行  $n-1$  次递归调用（这种分析只适用于一维空间。在三维空间中，为  $3+2(n-1)$  次调用。无论哪种情况，复杂度都为  $O(n)$ ）。

对于这些调用，组的平均大小  $m$  为  $n/2$ 。该函数最复杂的部分是排序，我们假定其复杂度为  $O(m \log_2 m)$  或  $O(n \log_2 n)$ 。这样，最糟情况下，总体复杂度为  $O(n) * O(n \log_2 n)$ ，即  $O(n^2 \log_2 n)$ 。

图 2.7.9 是时间复杂度接近最糟糕的情况（从物理上说，最糟糕的情况可能是不会发生的）。因为图 2.7.9 是一种不太可能出现的情况，因此，最糟糕的时间复杂度  $O(n^2 \log_2 n)$  有一定的误导性。有趣的是，在这种情况下，实际的时间复杂度为  $O(n^{1.78} \log_2 n)$ 。实际上，决不要期望最糟糕的情况会出现，因为这种情况极其罕见，且是人为的。

对于 RDC 而言，一种更可能发生的糟糕情况是所有物体都连接在一起。在这种情况下，沿每维进行分析时，得到的都是一组。这样，时间复杂度为  $O(3n \log_2 n)$ 。然后，整个组将被传递给蛮干比较函数，以找出相互碰撞的物体对。这实际上是最糟糕的情况，其时间复杂度为  $O(3n \log_2 n + (n^2 - n)/2)$ ，即  $O(n^2)$ 。由于该时间复杂度与蛮干比较算法相同，同时不可能存在同其他所有物体都相碰撞的物体，因此 RDC 算法的速度总是更快。



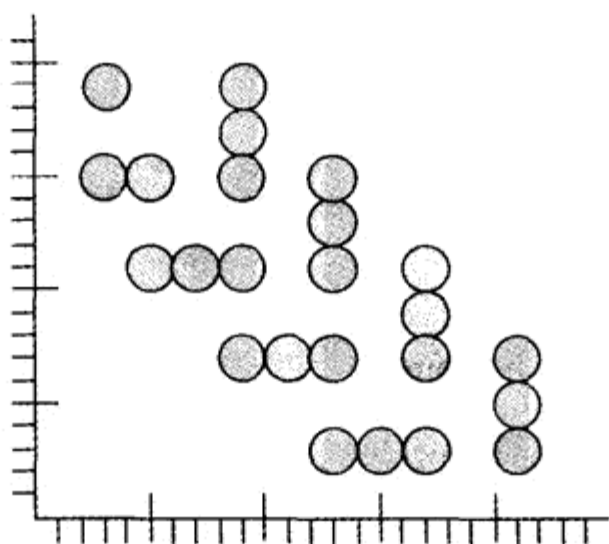


图 2.7.9 对于 RDC 而言，几乎最糟糕的情况（共 24 个物体）

### 2.7.5 结语

虽然递归逐维分组（RDC）算法并非一种完整的碰撞检测解决方案，但它是一个强大的过滤器，能够确定可能碰撞的物体组。其速度很快，远远超过了蛮干方法。虽然对空间进行分区是提高碰撞检测速度的标准解决方案，但 RDC 在不使用这种技术的情况下也获得了非常高的性能。然而，在实际应用中，RDC 最适合用于分析大量物体（超过 25 个）的碰撞，这可能应在使用分区算法减少可检测组后再进行。

### 2.7.6 参考文献

[Blow97] Blow, Jonathan "Practical Collision Detection," Proceedings, (Computer Game Developers Conference 1997), also available online at [http://142.104.104.232/eCOW/projects/Resources/practical\\_collision\\_detection.html](http://142.104.104.232/eCOW/projects/Resources/practical_collision_detection.html).

[Bobic00] Bobic, Nick, "Advanced Collision Detection Techniques," available online at [www.gamasutra.com/features/20000330/bobic\\_pfv.htm](http://www.gamasutra.com/features/20000330/bobic_pfv.htm), March 30, 2000.

[Roberts95] Roberts, Dave, "Collision Detection," *Dr. Dobb's Sourcebook* (May/June 1995): pp. 7~11.





## 2.8 不规则碎片编程

Jesse Laeuchli

jesse@laeuchli.com

当前，很多三维游戏和应用程序都显示随机的地形。无规则碎片（fractal）常被用来创建这样的地形。这种地形以高质量和逼真而著称，但只到最近计算机的处理能力才使得游戏机能够实时地处理这样的地形。

想使用无规则碎片来创建地形的原因很多。如果创建一个完全随机的地形，将山谷和山丘的高度设置为随机的，则地形看起来是随机的，但不真实。另一方面，如果使用一个有条理的函数（如正弦函数）来设置地形的高度，则地形将是有序且可预测的，这在现实生活中是不可能出现的。可以使用无规则碎片对这两种情况进行折衷，得到这样的地形：其中的山谷和山丘不是随机的，但更像我们看到的地形。

对于无规则碎片，一种不错的非数学定义是，以可控的方式加入了随机性的规则的东西。使用无规则碎片构建的东西并非完全随机的，但也不是完全有序的。无规则碎片最常被用作纹理（texture）或用于创建几何模型。不能使用无规则碎片来构建复杂但没有任何规律性的东西（如人）。

本文介绍一些比较有用的无规则碎片，并探讨如何使用它们。这里不打算探讨各种类型的无规则碎片，因为很多无规则碎片涉及高深的数学知识，而且在当前的游戏编程中的用途不大。在游戏中，无规则碎片通常被存储为高度图（height map）——一个矩形网格，其中每个单元格中都有一个表示该点的值的数字。就无规则碎片地形而言，这个值表示的是高度；对无规则碎片云而言，它表示的是密度。对这种值的解释方式有很多，其中的一些方式可能还没有发明，然而，由于最常见的情况是高度，因此作者将使用无规则碎片中某点的“高度”，而不是中性词“值”。

### 2.8.1 无规则碎片

最常见的一种无规则碎片是等离子体无规则碎片（plasma fractal）。这种无规则碎片很容易进行编程，作者对其前途最为乐观。读者可能见过数百个这样的无规则碎片（见图 2.8.1）。

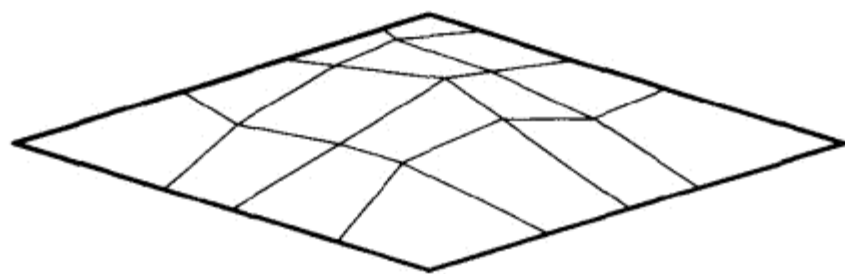


图 2.8.1 一个小型的等离子体无规则碎片

要创建等离子体无规则碎片，可以创建一个空的高程图，并指定4个角的高度。然后，对于每条边，根据两个角进行线性插值。四边形中心的高度为4个角高度的算术平均值。接下来，对这些新顶点定义的每个四边形重复上述过程。这样不断重复下去。这种插值方案将在4个角之间生成一个平滑的曲面；事实上它本身便可用于在高程图中进行平滑插值。

为使用这种简单的曲面来创建无规则碎片，我们在每次递归中都增加一个新的步骤。生成中央顶点后，将其高度增加或减去一个随机值。在下一次递归中，我们执行同样的操作，但将随机数的范围减去 $H$ 。注意，这是惟一控制无规则碎片的因素。如果 $H$ 为0，则高程图的变化很大； $H$ 越大，高程图的变化越小。

等离子体无规则碎片的主要优点是，它易于理解，运行速度也非常快。因此，除非速度的重要性高于结果，否则应避免使用无规则碎片，其原因有两个：首先，无规则碎片的可控程度很低；其次这种无规则碎片的人造痕迹非常明显——具体地说，就是山脊总是在边缘。

## 2.8.2 断层无规则碎片

另一种常见的无规则碎片是断层无规则碎片（fault fractal）。对它的编程也很容易，但速度比等离子体无规则碎片慢得多。断层无规则碎片模仿了强烈地震对随机“断层线”的影响。

要创建这种无规则碎片，同样可以先创建一个空的高程图，然后应用一系列的“地震”：创建一个穿越高程图的随机线，然后将随机线一边的每个单元格的值增加一个很小的值，将另一边的每个单元格的值减去一个很小的值。不断重复这一过程，直到对地形满意为止。

这种技术能够创建出一些非常好地无规则碎片，但创建速度非常慢。通常，要获得满意的结果，需要1 000~10 000条断层线。因此，无法实时地创建这样的无规则碎片！然而对于离线处理，这种无规则碎片实现起来非常容易，而且不会有诸如等离子体无规则碎片等细分无规则碎片的混层效应（aliasing effect）。

断层无规则碎片是为数不多的几种适用于球形的无规则碎片之一。在这种情况下，断层线是球上的大圆，将球分为两个半球。然后，每个半球在其公共平面上细微地移动。

除了使用的随机数序列外，断层无规则碎片没有控制参数，因此难以人为地生成。

## 2.8.3 FBM

前面介绍的方法都是专用的生成无规则碎片的方法。FBM（Fractal Brownian Motion）无规则碎片在数学方面更严密，并有非常良好的数学特性，因此使用起来更简单。

FBM无规则碎片是一系列噪声函数的组合，但其中的每个噪声函数都是独特的。这种技

术的关键是理解不同的噪声类型。

首先，有白噪声，它是完全随机的。电视上没有携带信号的静电噪声就属于白噪声。白噪声的高度变化很大，与附近点的高度没有任何关系。

其次，有粉红噪声（pink noise），其不同点的高度差在一定的范围之内。

在计算机图形学中，没有特别声明的情况下，术语噪声通常指的是粉红噪声。Ken Perlin 是在计算机图形学中使用粉红噪声的鼻祖，当时正编写现在非常著名的 noise() 函数（这种噪声被称为 Perlin Noise）。他也是第一个获得奥斯卡奖的图形程序员——以奖赏他编写了前面的函数（作者期望自己因以最佳的方式调用 fopen() 函数获得奥斯卡奖）。

虽然我们通常认为噪声是一维（如声音）或二维函数（如高程图），但它也可以是三维、四维，甚至更多维的。如果要在动画纹理（能够模拟火焰）或三维动画密度函数（可模拟汹涌的云）加入时间分量，则这很有用。

要创建粉红噪声，可创建一个网格（或维数更高的数组），并在每个单元格中存储一个随机值。然后根据这些值在曲面的每一点进行插值，从而创建无规则碎片。定义这种噪声函数的特性为频率——网格点之间的距离的倒数。频率越高，粉红噪声越接近于白噪声。

创建各种频率的粉红噪声后，创建 FBM 无规则碎片将非常容易，只需将噪声函数返回的高度混合起来即可。最简单的情况是将噪声函数相加，但也可以将它们相乘或采取其他的合并方式。

FBM 的参数比噪声多些，它除频率外，还有音频程（octave）、幅度和 H 参数。音频程变量设置将多少个噪声加到一起；幅度是调整噪声总体高度的变量；H 控制每隔一个音频程幅度变化多少。对每个噪声片段，都必须设置频率。为了解频率对 FBM 的影响，来看无规则碎片地形。将噪声加到一起时，频率低的噪声创建山丘和山脉；而频率高的噪声创建崎岖程度。这样，可以很容易地控制无规则碎片的外观。另外，也可以设置噪声函数不同部分的频率差，这样可以使无规则碎片的某些地方比其他地方更崎岖，从而得到山区较崎岖而沙漠较平坦的地形。

也可以将噪声函数相乘，而不是将它们相加。这样，无规则碎片将更为变化多端。将噪声相乘可以减略低的部分，增强高的部分。如果得到的高程图被用于地形，则地形中将有平原、山脉和丘陵——如果将噪声相加，则地形将有统一的外观。将噪声相乘时一定要小心，因为很容易走向极端，过度地减略和增强，结果得到一个包含一些尖峰的平原！

## 2.8.4 实现

探讨关于创建 FBM 的理论后，接下来实现一个噪声生成器。要生成噪声，必须生成一个随机数源。作者编写的随机数生成器的代码如下：

```
float random (int x , int y)
{
    int n=x+y*57;
    n=(n<<13)^n;
    float ret;
    ret= (1 - ( (n * (n * n * 19417 + 189851) + 4967243) & 4945007) / 3354521.0);
    return ret;
}
```

```
}
```

该函数将数字同质数相乘、相加和相减，并返回一个-1到1之间的伪随机值。当然，这并非生成随机数的惟一方式。有些系统使用预先生成的随机数，这样做对性能的提高不大，但必须存储大量的随机数。采用本文介绍的方法，只需记录用于重新生成随机数的种子即可。

接下来，需要一个对这些随机值进行插值的函数。一个简单的 `lerping` 函数就管用，其速度快，但提供的结果不是太好。虽然有些人使用样条线插值，但这是最慢、最复杂的方式。本文使用余弦差值来构建噪声函数，其质量高、速度快，且易于理解。注意，这并非惟一可用的噪声函数，它不同于 Perlin Noise，后者复杂得多。有关完整的列表，请参阅[Ebert98]。

下面是余弦差值的代码片段：

```
double cosineinterpolation(double number1,double number2,double x)
{
    double ft;
    double f;
    double ret;
    ft = x * 3.1415927;
    f = (1 - cos(ft)) * .5;
    ret=number1*(1-f) + number2*f;
    return ret;
}
```

有了随机数函数和插值函数后，便可以通过生成随机数并对其进行插值来创建噪声函数了。

```
float noise(float x, float y)
{
    int xinteger=x;
    float fractionx=x-xinteger;
    int yinteger=y;
    float fractiony=y-yinteger;
    float v1,v2,v3,v4,i1,i2;
    float ret;
    v1= randomnumber (xinteger, yinteger);
    v2= randomnumber (xinteger + 1, yinteger);
    v3= randomnumber (xinteger, yinteger + 1);
    v4= randomnumber (xinteger + 1, yinteger + 1);
    i1= cosineinterpolation (v1,v2,fractionx);
    i2= cosineinterpolation (v3,v4,fractionx);
    ret= cosineinterpolation (i1,i2,fractiony);
    return ret;
}
```

上述函数接受两个参数。通常，对于高程图或存储结果网格中的每个点，都应调用该函数。注意，这是一个二维噪声。

有些情况下，修匀噪声更合适。修匀将降低噪声的频率，使之看起来不那么方(square)。修匀一维噪声是没有意义的，因为通过降低频率可获得相同的效果。如果要修匀噪声，可在噪声函数中调用 `smoothrandom` 函数，而不是 `randomnumber` 函数。

```
float smoothrandom(int x,int y)
{
    float corners=(randomnumber(x-1,y-1)+
        randomnumber(x+1,y-1)+randomnumber(x-1,y+1)+randomnumber(x+1,y+1))/16;
    float sides  = (randomnumber(x-1, y)+randomnumber(x+1, y)+
        randomnumber(x, y-1)+randomnumber(x, y+1) ) / 8;
    float center = randomnumber(x, y) / 4;
    float ret=corners+sides+center;
    return ret;
}
```

这相当于对中央点周围的 9 个点进行积分。

建立噪声函数后，创建 FBM 将非常容易，其代码如下：

```
float FBM(float x, float y, float octaves, float amplitude, float frequency, float h)
{
    float ret=0;

    for(int i=0;i<(octaves-1);i++)
    {
        ret +=( noise (x* frequency, y* frequency)* amplitude);
        amplitude*=h;
    }
    return ret;
}
```

该噪声中的值随倍频程而变化的方式虽然在很多情况下都可行，但有时候对其进行修改可能会有帮助。您可以修改频率和幅度随倍频程的变化量，甚至可以忽略倍频程。这种控制是 FBM 的另一个优点（图 2.8.2）。

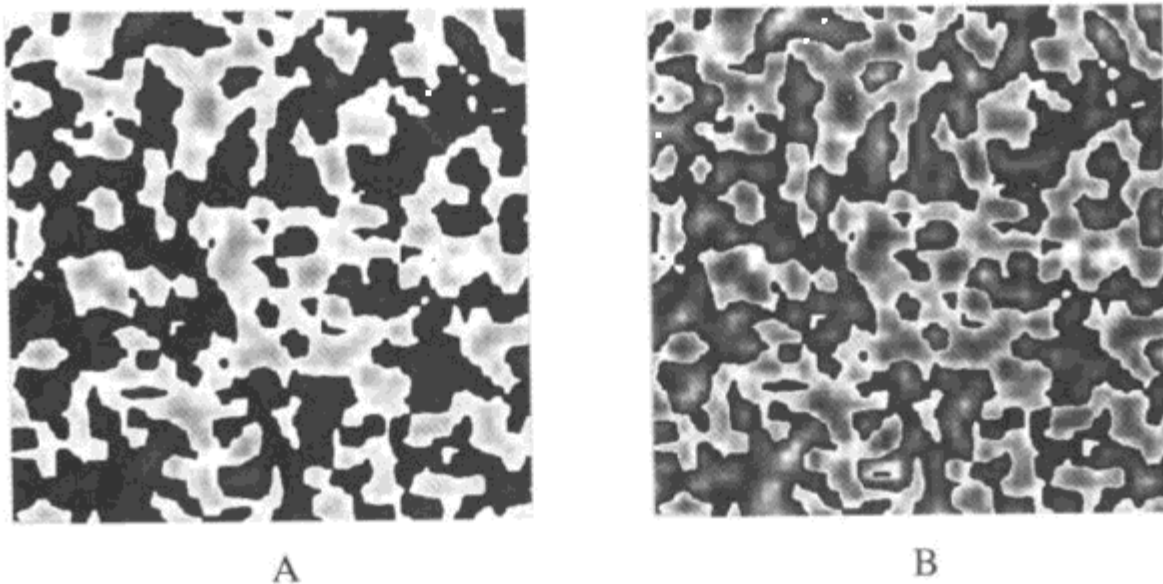


图 2.8.2 A)只有几个倍频程的 FBM; B)相同的 FBM 再增加几个倍频程后

FBM 的一种变体是多重无规则碎片 (multifractal)，它类似 FBM，只是不将噪声相加，而是相乘。下面是创建这种无规则碎片的代码：

```
float Multifractal(float x, float y, float octaves, float amplitude, float frequency,
```



```
float h, float offset)
{
    float ret=1;

    for(int i=0;i<(octaves-1);i++)
    {
        ret *=(offset)*( noise (x* frequency, y* frequency)* amplitude);
        amplitude*=h;
    }
    return ret;
}
```

`offset` 参数对噪声相乘进行一定程度的控制。您将发现，创建的无规则碎片中的丘陵、山脉和平原完全不同。

### 2.8.5 使用 FBM

下面介绍如何使用 FBM 来创建云。

#### 1. 云

使用 FBM 来创建云非常容易。首先生成一个 FBM，然后将高度大于 0 的点视为白色，并给它指定一个大于 0 的 `alpha` 值；对于所有小于 0 的高度，将其 `alpha` 值指定为 0。将该纹理贴到四边形或球形上，然后根据建模时的天空情形进行着色。图 2.8.3 是一些云纹理，详细的信息请参阅 `clouds.cpp`。

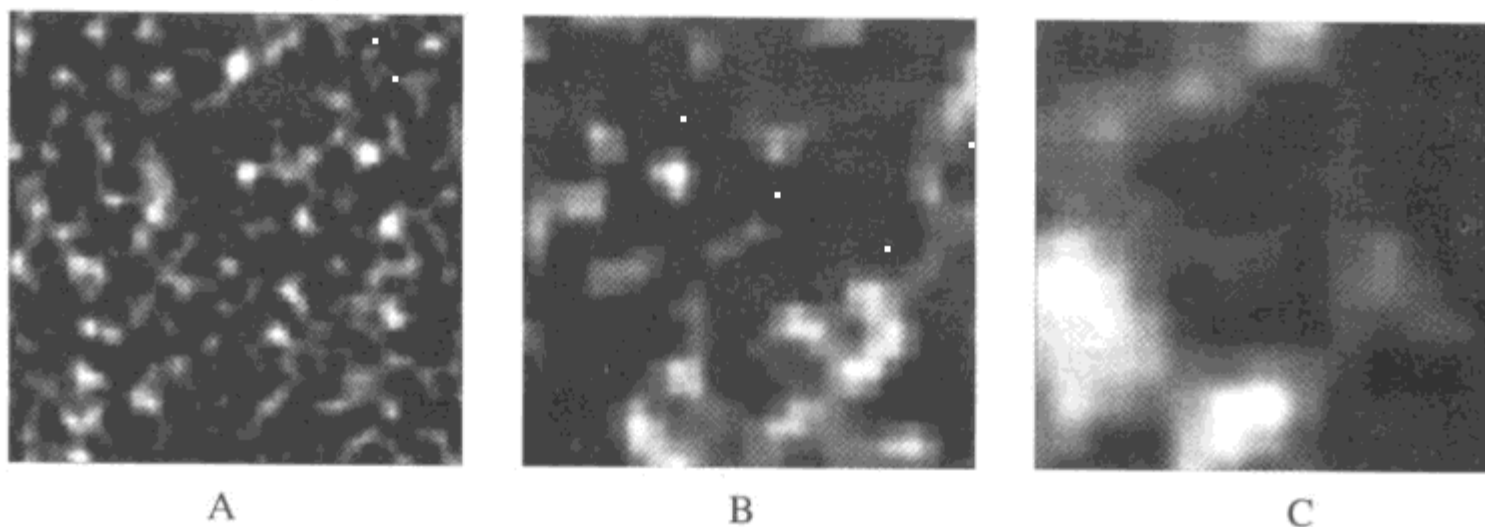


图 2.8.3 使用 FBM 生成的云（降低了频率）

#### 2. 地形

对创建地形，FBM 也很有用；事实上，FBM 可用于创建卓越的地形。只需将高程图解释为地形的高度（高程图因此而得名）即可。要显示地形，可使用三角形创建一个平面，在三角形的每个顶点上，使用高程图中的高度（图 2.8.4 和 2.8.5）。详细的程序清单，请参阅 `landscape2.cpp`。

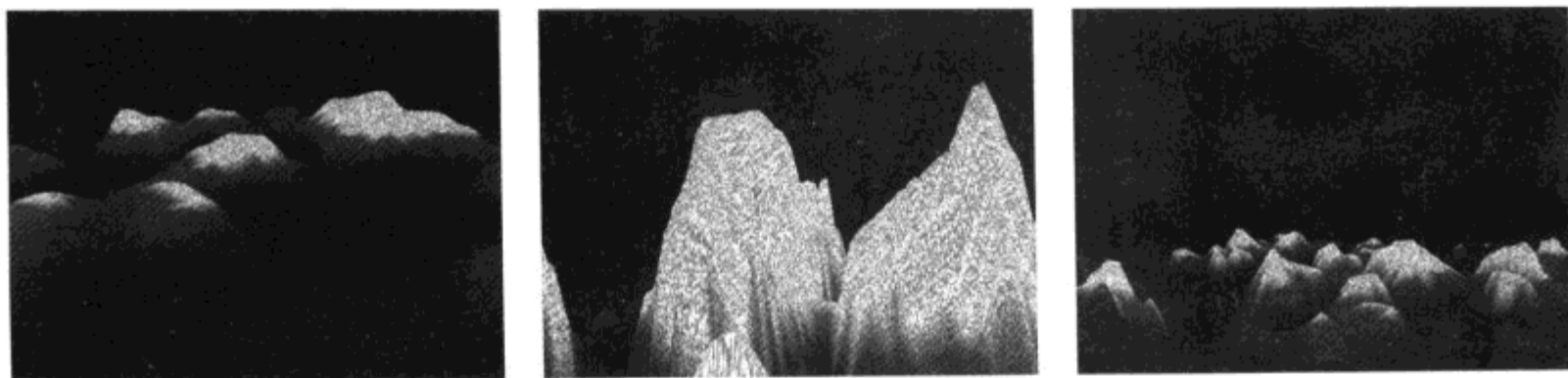


图 2.8.4 各种无规则碎片高程图：地面纹理是一个按高度着色的高频率 FBM。

A) 使用低频率、低倍频程和小幅度生成的；B) 幅度更大时； C) 频率更高时

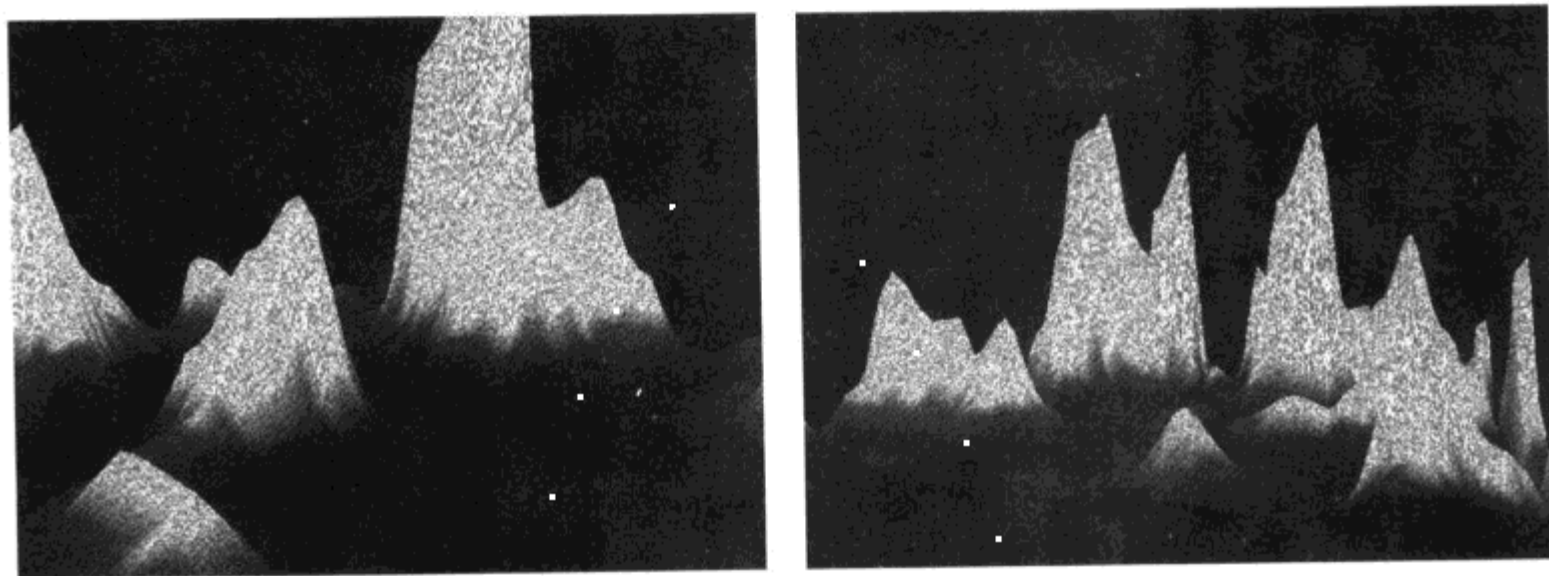


图 2.8.5 两个多重无规则碎片地形。注意，同一个地形中包含山谷、平原、丘陵和高山

## 2.8.6 参考文献

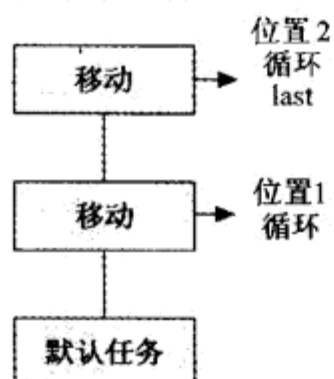
[Ebert98]. David S. Ebert, et al. *Texturing and Modeling*. San Diego: Academic Press, 1998.



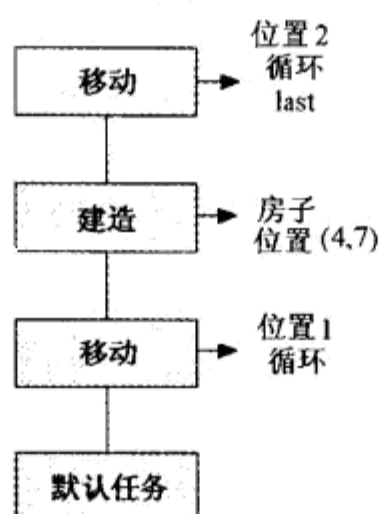


# 人工智能

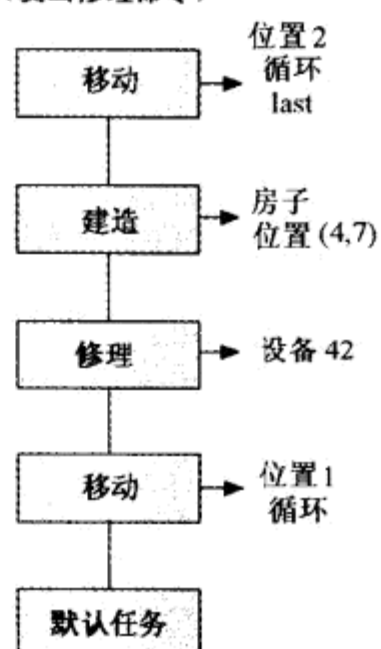
大脑队列  
(发出巡逻命令)



大脑队列  
(发出建造命令)



大脑队列  
(发出修理命令)





绪 论

Steve Rabin, 任天堂 (美国) 公司  
steve\_rabin@hotmail.com

欢迎阅读《游戏编程精粹 2》的“人工智能”一章。《游戏编程精粹 1》取得了极大的成功，尤其是其中关于人工智能的一章，因此本章任重道远。所幸的是，接下来的所有文章都是创新性的，是各位作者经验、智慧和努力的结晶。这些都是实用的 AI 技术，可被立刻用于几乎任何正在开发的游戏。

虽然本章中的文章是由 11 位不同的作者分别撰写的，但每篇文章都不是孤立的。其中的很多文章是相互依赖的，而其他的几篇文章则对《游戏编程精粹 1》中介绍的概念和技术进行了更深入的讨论。智慧的火花贯穿整章，其中每篇文章都值得阅读。图 3.0.1 说明了本章的组织结构和布局。

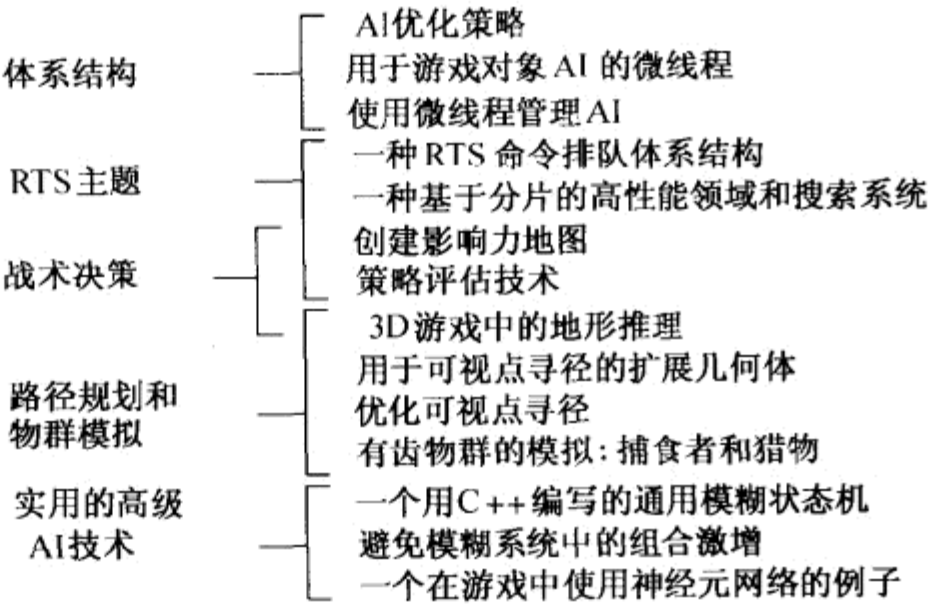


图 3.0.1 本章的组织结构

多年以来，游戏领域对来自学术研究的高级 AI 成果进行了激动人心的尝试，如遗传算法、神经网络 (Neural network) 和其他异乎寻常的技术。从那时起，钟摆便摆向了另一个方向，现在大多数 AI 游戏程序员都坚持使用更简单、经过实践检验的、总是管用的技术。

有趣的是，专注于经过实践检验的方法促使 AI 游戏程序员开发了一个独立于学术研究的领域。这没有什么好奇怪的，因为学术研究人员和游戏 AI 开发人员的目标完全不同。学术研究追求的是过程，关心如何解决问题；而游戏编程追求的是结果，关心最终产品的表现。因此，游戏编程

通常是一门不精确的科学，为实现游戏目的而在计算方面进行折衷；而学术研究决不追求折衷的解决方案，因为它真正重视的是算法。这两种做法都是正确的，因为它们的目标和动机不同。

A\*搜索算法是这种差别的典范。在游戏中，这种算法常被用来找出角色通过复杂环境的路径。学术界之所以对 A\*算法感兴趣是由于它能够找到最短路径，且计算速度比较快。然而，只需对其中的推理做细微的修改，便可极大地提高 A\*算法的速度。不幸的是，这样的修改将导致找出的路径不一定是最短的，不过这样的路径还是很不错的。

由于不能保证路径是最短的，学术研究人员认为这种修改后的推理是无法接受的，他们对这种近似但不是最佳的结果嗤之以鼻。因此，学术界的研究几乎不使用不可接受的推理。然而，在游戏中，这种无法接受的推理正是您想使用的，因为其速度快得多，而结果相差无几。有关 A\*算法的详细讨论，请参阅《游戏编程精粹 1》中的“A\*算法速度优化”一文。

正如您看到的，要推动这一领域向前发展，有待于游戏程序员的努力，有待于他们对其感兴趣主题进行改进。如果学术研究有助于在游戏中创建 AI 系统当然很好，但游戏程序员也需要进行自己的研究。因此很明显，游戏开发领域需要像本书这样的图书。从这种意义上说，探索和促进这样的技术并进行有用的游戏开发研究就是本书的职责所在。



## 3.1 AI 优化策略

Steve Rabin, Nintendo of America

steve\_rabin@hotmail.com

**复**杂的 AI 需要占用大量的计算时间。当数以十（百）计的自主体（autonomous agent）必须同时在场景中智能地漫步时，情况将非常糟糕。然而，这并非通常的优化问题。游戏的很多领域都涉及到可伸缩性，而 AI 还支持虚拟的并行性。这种并行性源自每个 AI 主体都运行自己的代码，因此好像所有的主体都在同时思考一样。

从其出现起，并行性就是 AI 可供利用的重要资产之一。可以通过操纵 AI 主体的并行性和其他独特的特性来优化大多数 AI 系统。阅读下述优化策略时，别忘了有数百个 AI 主体，它们的行为和日程各不相同。在这种情况下，一些策略的合理性将更为明显。

### 3.1.1 策略 1：使用事件驱动行为而非轮询

最理想的情况是，自主体应不断地监控其所处的环境，并做出相应的反应。在每一帧中，它们都必须对环境进行轮询，以获悉自己需要对哪些事件或对象进行响应。这意味着从碰撞检测（collision detection）到注意感兴趣的事情等各项内容。这种方法的问题在于，它将导致大量重复计算。通常，分别轮询是不经济的。

替代方法是尽可能采用事件驱动技术。例如，在棒球游戏中，当棒球被击打时，棒球只需告诉接球手，自己被击打即可。与每个接球手对环境进行轮询，以了解击球手何时击打了棒球相比，这样做的效率要高得多。当然，这显然是一种事件驱动行为，但它表明，通过使用事件驱动策略，可以极大地节省处理时间。

另一个例子是战场上飞行的箭。通常，箭将被检测是否同角色（character）发生碰撞。如果箭射中了某个角色，它将指出该角色及其被射中的部位。然后，该角色将根据被射中的部位做出相应的反应。如果不采用事件驱动方式，这将是引人入胜的。

现在假设您希望角色对射来的箭做出反应，甚至试图避开之。如果采用轮询解决方案，则每个角色必须间歇性地考虑附近的箭。相反，如果采用事件驱动解决方案，则可以让箭的碰撞检测函数对将发生的碰撞进行预测，并指出将可能被射中的角色。角色可以根据信息做出任意反应，包括逃跑、蹲下或呈现恐惧的表情。这样行为将是丰富多彩的，而无需对每种

意外情况进行轮询。

在《游戏编程精粹 1》的“设计一个通用、健壮的 AI 引擎”[Rabin00a]一文提供了一个典型的事件驱动 AI 系统。

### 3.1.2 策略 2: 减少重复计算

---

这种策略旨在通过在多个 AI 主体之间共享结果以减少重复计算。很多时候, 虽然本可以只计算一次并共享计算结果, 但主体将重复计算相同的信息。这种优化很容易, 通常可以节省大量的时钟周期。

这种策略的一个简单范例出现在碰撞检测 (collision detection) 中。如果每个物体都各自进行碰撞检测, 则检测次数为  $n^2-n$ , 也就是说, 如果有 100 个物体, 则需要进行 9900 次碰撞检测。然而, 由于对于给定的每两个物体, 只需进行一次碰撞检测即可, 因此优化的碰撞检测算法只需进行  $(n^2-n)/2$  次碰撞检测, 也就是说, 如果有 100 个物体, 则只需进行 4950 次碰撞检测。这种节省直接源自于消除了冗余计算。

另一个消除冗余计算的例子可在寻径 (pathfinding) 中找到。当玩家命令一系列的个体 (unit) 都移到地图的另一侧时, 每个个体可能各自进行寻径。而一种更简单、更快的解决方案是, 让其中的一个个体进行寻径, 然后让其他个体跟在他的后面。这样可以避免几乎相同的路径请求, 从而节省大量的时钟周期。

### 3.1.3 策略 3: 由管理员集中进行协调

---

主体常常需要同其他主体进行合作, 无论它们是突击小分队还是体育小组。对于多个主体之间的协作, 可由一个管理员实体来做出复杂决策, 从而使协作更简单、更快速。这种复杂的决策通常是确定每个成员的任务 (role), 而主体将完成相应的任务。

例如, 假设有一个潜入到一幢大楼中的小组。为确保每个位置的安全并向前移动, 小组成员必须紧密协作。如果成员必须各自决定其行为, 则协调起来将非常困难, 同时成员之间必须进行大量的交流。相反, 可由管理员实体对行动的每一步进行策划, 并指挥各个成员接下来应采取的行动。这种方法更简单, 使设计更健壮、更有效。[Rabin98]中提供了一个管理员如何简化棒球游戏的范例。

使用这种策略时, 无需在屏幕上将管理员显示出来, 这至关重要。这些管理员通常是一个虚构的实体, 他们对 AI 主体之间复杂的协作进行组织。您甚至可以虚构临时的管理员, 他们动态地组织一系列团结在一起的怪物小组。如果接下来该小组解散或与另一个小组合并, 则可以给新的小组指定管理员。通过这些管理员, 各个怪物将协调一致, 选择最佳的目标, 而不是向同一个敌人发起进攻。

### 3.1.4 策略 4: 不那么频繁地运行 AI

---

AI 主体很少需要在每帧都运行其整个决策函数。很多时候, 主体可以每隔几帧或几秒运行代码的一部分。因为真正的生物都有一定的反应时间, 因此 AI 主体无需像闪电一样快速

反应。这样，便提供了一种简捷的削减 AI 处理的途径。

一种实现这种策略的方式是，使用支持任意定时器回调（timer callback）的主体体系结构。如果主体能够轻松地设置定时器，并在定时器超时后通知主体，则可以构建一个灵活、可调的系统。《游戏编程精粹 1》中的“设计一个通用、健壮的 AI 引擎”一文[Rabin00a]讨论了一个定时器消息接发系统，非常适合用于实现这种策略。

这种策略的问题之一是，可能存在 AI 处理高峰。如果大多数主体同步执行其回调、同时每隔  $N$  秒执行一次，将发生这样的问题。对此，最简单的解决方案是，随机化每个主体的处理窗口。例如，主体可以每隔 0.3~0.5 s 执行其入侵检测代码，并在每次回调后，随机选择新的延迟。这种随机执行窗口实际上确保了主体之间不会同步，从而导致处理高峰产生。

### 3.1.5 策略 5：将处理工作分散到多帧中完成

---

通常，A\*寻径是非常可怕的算法之一，将占用大量的 CPU 周期。由于在几帧之内，情况变化不会太大，因此主体可以将寻径计算工作分散在几帧中完成。通过保存每帧中计算得到的结果，并在下一帧恢复该结果，可在 2~4 帧之内找到路径。这样，每帧的处理负载将更低。任何所需时间不确定的算法，都可按这样的方式进行分解。

这种策略可作为一个特例在一个模块（如寻径模块）中实现，也可作为 AI 操作系统的一部分。本章的下一篇文章介绍了如何实现微线程（micro-thread），以最大限度地降低 AI 负载；这种体系结构使得 AI 主体能够更容易地将计算分散到多帧中完成，并有助于控制处理负载。再下面的一篇文章描述了如何在微线程环境中以最佳的方式组织 AI。

### 3.1.6 策略 6：利用细节级 AI

---

细节级（level-of-detail, LOD）概念是一种精巧的优化，被大量用于图形领域之外。当前，大多数游戏引擎都使用图形 LOD，以便快速显示远处的物体。其思想是，当物体远离摄像机时，使用一个包含较少多边形的模型；而当物体较近时，使用一个包含较多多边形的模型。这样，显示速度将极大地提高，而对游戏的视觉效果几乎没有任何影响。需要指出的是，细节级概念也可用于其他游戏编程领域，如 AI。

实际上，用于 AI 的细节级技术可归结为三种策略：首先，根据主体离摄像机、玩家或现场的距离，改变其处理频率；其次，根据相关程度，改变主体算法的复杂性，方法是如果主体远离画面，则不进行精确的寻径；第三，当主体对玩家的重要性降低时，在单个模拟算法中显示多个主体。对于第三种变体，一个极端的例子是，使用一个简单的公式来模拟远处争斗的结果，而对于离玩家较近的战斗，则对每个士兵和每颗子弹都进行模拟。

细节级旨在在玩家不会觉察的情况下减少计算量。如果玩家将发现问题或有差别，则说明优化过头了，就像在视觉上出现跳跃的图形 LOD 一样。

### 3.1.7 策略 7：只解决问题的一部分

---

对于大型问题，有时候只需解决亟需解决的部分便足够了，可在以后需要时再计算解决



方案的其余部分。甚至可能出现这样的情况，情形发生了很大的变化，问题的余下部分已经无关紧要，不再需要计算。这种策略非常类似于迟缓计算（lazy evaluation）。

这种策略的一个典范是层次式寻径。在层次式寻径中，路径寻找分为两步。首先，计算角色到达目的地的高级（从房间到房间）路径；然后，进入到每个新房间后，再计算让角色到达目的地的下一个房间的微路径。如果由于某种原因，角色改道了，则丢弃余下的路径，这些路径不会被计算。这种按需寻径可以极大地提高速度，对于需要在大型区域中导航的游戏来说，这至关重要。有关层次式寻径问题及其实现的完整描述，请参阅[Rabin00b]。

### 3.1.8 策略 8：离线完成困难的工作

---

有时候，问题非常困难，没有解决它所需的 CPU 时间。在游戏开发的早期，诸如计算余弦的问题都将占用过多的时间，因此程序员们会预先计算某一范围的角度余弦，并通过一个包含索引的查找表来获得所需的值。虽然现在这种技术已不再重要，但其基本策略还是很有用的。

当前，这种策略的具体范例包括预先计算 BSP、预先分析地形以及经过仔细训练的神经元网络。后面有三篇举例说明这种策略的文章：William van der Sterren 撰写的“3D 动作游戏中的地形推理（terrain reasoning）”、Thomas Young 撰写的“可视点寻径中的扩展几何体”和 John Manslow 撰写的“在游戏中使用神经元网络：一个具体的范例”。这些文章都探讨了大量可用于分析、精炼和存储专用信息的离线时间，在运行期间可以使用这些信息。这无疑是最强大的优化策略之一，因为它能够将几千小时的计算结果存储为几 KB 的数据，而装载这些数据只需几个 CPU 周期。

### 3.1.9 策略 9：使用突发行为以避免编写脚本

---

这一策略与前一策略相反。您常常没有足够的时间来为出现在场景中的数以百计的背景 AI 实体创建行为或脚本。对此，解决方案是制定非常简单的规则，以实现智能、突发的行为（emergent behavior）。本书后面的“有齿物群的模拟：捕食者和猎物”一文（由 Steven 撰写）详细介绍了物群模拟（flocking）技术如何帮助创建神奇、复杂的行为，这些行为无法通过编写脚本来完成，但有趣、逼真。

不幸的是，这种策略是一把双刃剑，在减少离线工作量的同时，可能导致意外的结果。Steven 描述了一个充满生物的世界，其中的一些生物依靠捕食其他生物来生存。食肉动物可能将猎物消灭殆尽，导致整个生态系统崩溃。从其本质上说，突发行为是不可预测的，难以进行详尽的测试。因此，它最适合用于非关键 AI 实体，这些实体并非游戏发展或完成的关键，如背景中的野生生物。

### 3.1.10 策略 10：通过连续记录来分摊查询成本

---

有时候，需要收集大量的数据，以做出智能的决策。如果需要数据时才进行收集，则用于计算的时间可能是不可接受的。对此，解决方案是不断地更新数据结构中的数据。虽然跟

踪数据所需的时间和内存更多，但这被分摊到很多帧中。这样，对数据的简单查询将不会影响帧频。

影响力地图 (influence map) 是这种策略的一个典范。影响力地图用于分析总体战略模式，以加快 AI 决策的速度。随着游戏的进展，个体 (unit) 在移动、演化或消失时，将在影响力地图中更新其信息。这样，当总体战略 AI 需要分析场景时，数据是现成的。因此，能够快速查询影响力地图，而不会对速度有太大的影响。本书后面的“构建影响力地图”一文 (由 Paul Tozour 撰写) 解释了何为影响力地图，并就如何使之更强大提出了一些独特的见解。

这种策略也是本书后面的“一个基于分片地图的高性能视域 (line-of-sight) 和搜索系统”一文的关键所在。在该文中，Matt Pritchard 使用连续记录的策略来实时地维护数百个主体的视域 (line-of-sight) 数据，这并非一项简单的任务。很多 RTS 游戏在战争迷雾 (fog-of-war)、袖珍地图和总体智能方面，并未实现帧级 (frame-by-frame) 准确性；然而 Matt 完整地解释了使《帝国时代》系列在这方面出类拔萃的技术和策略。

### 3.1.11 策略 11：重新考虑问题

---

Michael Abrash 是游戏开发领域的优化大师，同时是图形研究领域的领军人物之一。他多次强调 (通过演讲和撰文 [Abrash00]) 重新考虑问题以将代码的速度提高多个数量级的重要性。其主要观点是，优化代码的某部分常常会获得意外的收获。真正的优化方法是，从稍微不同的角度，以其他方法或算法来解决问题。

本书后面的“避免模糊系统中的组合激增 (combinatorial explosion)”一文 (由 Michael Zaroinski 撰写) 解释了这种策略。在该文中，Michael 解释了一种名为 Combs Method 的替代算法，这种算法巧妙地避免了传统模糊系统的指数特征。虽然这种算法得到的结果不相同，但具有可比性，对大多数用途而言足够了。另外，它简化了模糊逻辑的实现，让任何人都能够轻松地将其集成到游戏中。

从不同的角度重新考虑问题，可能是任何人都能够提供的优化建议。它是其他优化的根源所在。只有重新定义或抽象问题、进行类比或改变您的角度，才能完成真正优化代码的一跃。虽然这种跳跃并非常常出现，但通过阅读和学习其他人如何解决类似的问题，可以从他们的跳跃中获得灵感。

### 3.1.12 结论

---

优化 AI 系统面临的问题时，应从稍微不同的角度进行思考。当您下次试图提高一个看起来不可优化的系统的速度时，别忘了阅读上述策略清单。通过从许多不同的角度考虑问题，该清单将有助于您温故知新，并思考如何将每种策略应用于具体的情况中。下面再次列出了这些策略：

- (1) 使用事件驱动行为而不是轮询；
- (2) 减少冗余计算；
- (3) 由管理员集中进行协调；
- (4) 不那么频繁地运行 AI；

- (5) 将处理工作分散到多帧中完成;
- (6) 利用细节级 AI;
- (7) 只解决问题的一部分;
- (8) 离线完成困难的工作;
- (9) 使用突发行为以避免编写脚本;
- (10) 通过连续记录来分摊查询成本;
- (11) 重新考虑问题。

阅读后面的文章时, 请考虑其中应用了哪些优化策略。存在各种各样的策略, 令人惊讶的是不同的问题常常可以通过相同的策略来解决。只有天才才能洞察其中的联系。

### 3.1.13 参考文献

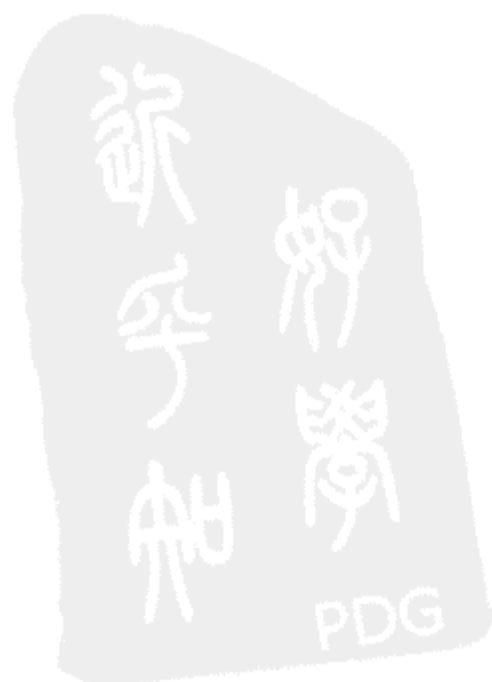
---

[Abrash00] Abrash, Michael, "It's Great to Be Back! Fast Code, Game Programming, and Other Thoughts from 20 (Minus Two) Years in the Trenches," Conference Proceedings, (Game Developers Conference 2000). Text available online at [www.gamasutra.com/features/20010117/abrash\\_01.htm](http://www.gamasutra.com/features/20010117/abrash_01.htm). Video also available online at [www.gamasutra.com/features/index\\_video.htm](http://www.gamasutra.com/features/index_video.htm).

[Rabin98] Rabin, Steve, "Making the Play: Team Cooperation in Microsoft Baseball 3D," Conference Proceedings, (Computer Game Developers Conference 1998).

[Rabin00a] Rabin, Steve, "Designing a General Robust AI Engine," *Game Programming Gems*, Charles River Media, 2000: pp. 221~236.

[Rabin00b] Rabin, Steve, "A\* Aesthetic Optimizations," and "A\* Speed Optimizations," *Game Programming Gems*, Charles River Media, 2000.



## 3.2 用于游戏对象 AI 的微线程

Bruce Dawson, Humongous Entertainment

bruced@humongous.com

**编**写优秀的 AI 非常困难。花费多年的时间和数百万美元后，计算机程序才击败了世界上最优秀的国际象棋选手。游戏 AI 不一定要参加国际象棋锦标赛，但必须每帧更新大量的游戏对象，且花费很少的 CPU 时间。

要编写优秀的 AI，除了其本质决定的复杂性外，常用的实现游戏 AI 的方法还将引入“偶然复杂性”（accidental complexity）[Brooks95]。

作为大多数游戏代码中的一个偶然复杂性的例子，假设我们有一个“看门人”对象，它在场景中游弋，做清扫工作。该看门人的例程很简单：选择目标并向它移动，当离目标足够近后便处置之，并不断重复这种过程。当然，由于这是游戏，我们希望这样的工作分小步完成，每个更新循环完成一步。其 C++ 伪代码（没有偶然复杂性）可如下：

```
void Janitor::Process() {
    while (true) {
        GameObject* target = GetNewTarget(this);
        while (Distance(this, target) > k_collisionTolerance) {
            WaitOneFrame();
            MoveABitTowards(this, target);
        }
        Dispose(this, target);
    }
}
```



然而在游戏中，上述代码并不太好，因为它不能与场景中的其他实体共享 CPU。因此，实现这种看门人对象的传统方式是，编写一个 Janitor 类，其中包含一个在每帧中都被调用的虚函数，如下述代码（摘自本书附带光盘中的范例）所示：

```
Janitor::Janitor()
    : m_target(0), m_state(k_NeedsGoal) {
}

void Janitor::ProcessFrame() {
    switch (m_state) {
        case k_NeedsGoal:
            m_target = GetNewTarget(this);
    }
}
```

```

        if (Distance(this, m_target) <= k_collisionTolerance) {
            m_state = k_DisposeObject;
            ProcessFrame();    // Call ourselves.
            return;
        }
        m_state = k_MovingTowardsTarget;
        // Intentionally missing break.
    case k_MovingTowardsTarget:
        MoveABitTowards(this, m_target);
        if (Distance(this, m_target) > k_collisionTolerance)
            break;
        else
            m_state = k_DisposeObject;
            // Intentionally missing break.
    case k_DisposeObject:
        Dispose(this, m_target);
        m_state = k_NeedsGoal;
        break;
    }
}

```

太乱了！代码不再简单、精致，而成了一个庞大、难读的状态机。修改代码的第二个版本显然要复杂得多，同时该状态机不能被其他类重用，因为它需要 `m_state` 和 `m_target`。这就是偶然（非本质）复杂性。

问题的核心在于，现在必须在对象中存储所有的状态。这是个深奥的变化。在代码的第一个版本中，`target` 变量被声明为局部变量，仅在需要时被创建，并在其工作完成后消失。而现在，必须在对象中记录目标，同时来在 `Janitor` 类中添加了另一个成员变量 `m_state`。该变量来自何方呢？为何代码的第一个版本不需要它？

在 `Process()` 的第一个版本中，当前执行的代码行指出了状态，它被存储在程序计数器中。程序计数器跟踪当前状态，栈指针指向当前的变量组，因此很容易。当我们使用回调来实现 AI 时，必须模拟指令指针和栈指针以及它们带来的好处。

### 3.2.1 一个更简单的方法

至此，将回调用于实体代码将导致混乱这一点是显而易见的。那么，如何对看门人代码的简单版本进行修改，使之管用呢？为此，需要让游戏循环的其他部分和其他对象也能使用 CPU，同时需要同步化对象，使之每帧进行一次更新。为此，需要将每个游戏实体放在一个独立的线程中，并让函数 `WaitOneFrame()` 切换到下一个线程。有了这样的系统后，第一个范例的伪代码便能够被编译并完美地运行！

对于每个对象，启动一个 OS 线程。这样，每个线程都像是拥有 CPU，而操作系统管理改变 CPU 状态的工作。函数 `WaitOneFrame()` 将完成切换到下一个线程所需做的工作——有关细节请参阅范例代码。

对于 Web 服务器和其他多线程（Multithreading）应用程序而言，这样的策略工作得很好；但对于游戏来说，这却是一种糟糕的选择。很多游戏运行在这样的环境中，即其中的操作系



统不支持线程，甚至根本就没有操作系统。即使是支持线程，线程对游戏而言也过于昂贵。在 Win32 中，切换线程需要花费数以千计的机器周期，而每个线程至少耗用 8KB 的内存——其中 4KB 用于线程信息块（Thread Information Block），4KB 用于最小的栈。这两种内存分配（Memory allocation）都为 4KB，是因为为确保线程的独立性并允许栈增长，分配的内存被放在独立的内存管理页中，这使得粒度为 4K。

在 Win32 上，也可以使用 fiber——协作式多任务线程。就我们的目的而言 fiber 更好些，因为场镜（context）切换的速度快得多，也更容易控制。然而，栈仍然至少占用 4KB，而且在 Win95 中 fiber 根本不管用。在 Win98 中，每个 fiber 栈至少为 8KB，这使得内存使用情况与常规线程一样糟糕。

### 3.2.2 微线程

我们回过头来，仔细考虑我们的需求。我们希望能够编写 AI 或其他对象更新代码，使这些代码好像拥有 CPU。在一个时隙更新中执行这些代码后，我们希望能够调用 WaitOneFrame()，给游戏的其他部分一些 CPU 时间。WaitOneFrame() 返回后，我们希望在接下来的一个时隙中接着执行后面的代码。我们希望这种过程快速完成，并希望其内存开销尽可能小。

必须从一个线程切换到另一个线程吗？CPU 中的指令指针是一个寄存器，它指向接下来要执行的指令。如果编写一些修改指令指针的代码，则可以轻松地从一个线程跳到另一个线程。修改指令指针很容易——CPU 有很多完成这项工作的指令。只需几行汇编语言代码，便可获得当前的指令指针，将其存储到某个地方，然后跳转到新的位置。这可以实现我们期望的效果——跳回到以前运行的代码片段中，然而，如果仅仅这样做，您将极其失望。

指针指令并非 CPU 中惟一的状态信息，也不是决定将执行哪个函数的惟一因素。另一项重要因素是栈指针。函数中的所有局部变量都是相对于栈指针（或栈帧指针，我们暂时忽略这一点）来存储的。事实上，栈指针甚至比指令指针更重要，因为函数的返回地址——指令指针——被存储在栈上。

假设我们编写了一段这样的代码，它将指令指针压入到栈中，然后修改栈指针，然后从栈中弹出指令指针。由于新的指令指针是从新栈中弹出的，因此修改栈指针也就修改了指令指针，也就是说，改变了线程。事实上，由于调用函数时，将把指令指针压入到栈中，而从函数返回时，将弹出指令指针，因此我们的线程切换函数包含两条指令——将一个新值移动到栈指针中，然后返回！



然而，事情并没有那么简单。首先，除了栈指针和指令指针外，CPU 中还有其他寄存器。没有问题，我们将以同样的方式处理它们。开始时将它们都压入到栈中，修改栈指针，然后在结束时将它们都弹出。这样便可以了。具体的细节随 CPU 而异，但任何在指向栈指针方面足够灵活的 CPU 都能实现这样的系统。在任天堂的 GameBoy 上已经实现了微线程！在 X86 CPU 上，需要 12 条汇编指令便可实现完整的线程切换。在本书附带光盘中的微线程范例程序中，就只使用 12 条汇编语言指令。范例函数 SwapThreads() 的代码如下：

### SwapThreads

```
// Preserve all of the registers that VC++ insists we preserve.
// VC++ does not care if we preserve eax, ecx and edx.
push    ebx
push    ebp
push    esi
push    edi

// Swap stack pointers
mov     eax, esp
mov     esp, s_globalStackPointer
mov     s_globalStackPointer, eax

// Restore all of the registers that we previously preserved.
// Yes, they're coming off of a different stack - they were
// carefully placed there earlier.
pop     edi
pop     esi
pop     ebp
pop     ebx
ret
```

只需要这样做吗？不一定。在有些 CPU 上，可能还必须保护浮点寄存器或多媒体寄存器。在 Intel 体系结构中，由于浮点寄存器的奇特排列，编译器无法在函数调用时保护这些寄存器，因此编译器总是在调用任何函数（包括 `WaitOneFrame()` 函数）前确保所有的浮点数都被写入到内存中。也不必保护所有的整数寄存器，因为编译器不会期望调用函数时保护这些寄存器——有关细节请参阅编译器文档。

### 3.2.3 栈管理

然而，改变线程时如何指定栈指针呢？启动 OS 线程时，隐式地分配和初始化了一个栈——现在，我们显式地完成这项工作。在 Win32 上，操作系统为每个栈分配一个 4KB 的存储页面，并预留一定的地址空间（默认为 1MB），以便栈可以增大。如果在函数中声明了一个大型数组或进行了很深的递归函数调用，则操作系统将自动分配更多的 4KB 存储页面。如果超过了预留的地址空间，操作系统将通过页面错误来阻止。

我们要避免给每个线程分配 4KB 的栈，那么应分配多少内存呢？我们可能决定为每个对象分配少于 100 字节的栈，毕竟我们需要成百上千甚至成千上万个这样的东西。我们可以使用 `malloc` 或 `new` 来分配 100 字节的内存块，然后通过仔细地初始化设置该栈，以便能够使用 `SwapThreads()` 切换到该栈。这种方法管用，但很危险。如果您编写的游戏实体代码在使用栈时超越了其边界，则将带来非常严重的后果。这将破坏该线程块前面的内存块，导致游戏崩溃。如果决定使用这种系统，一定要在栈的末尾加上某种标记，并在每次线程切换后检查这些标记是否被覆盖。这样，当您破坏了内存中的数据时，将会知道这一点。

另一种稍微不同的微线程实现可避免这种严格的栈长限制。这种实现分配一个大型栈，供所有的微线程共享。当微线程进入休眠状态时，线程管理器将栈的内容复制到一个大小可



调的备份缓冲区中。缓冲区仅在线程的栈长增大时才需重新分配，因此缓冲区分配时间可忽略不计。所有，这种方法惟一新增的开销是来回复制栈的内容。有趣的是，这种复制实际上是无开销的，因为它预先准备好了 CPU 高速缓存，使之为运行线程做好准备。

至此，相对于为每个线程分配定长的栈，栈复制好像没有任何优势可言。然而，栈复制的优点在于，切换线程时只需少量地使用栈。如果您的 AI 实体需要调用复杂的 BSP 寻径例程、调试打印函数或其他大量使用栈的函数，则可以使用栈复制微线程。临时使用大型栈不会有什么害处，只要不在这些实用函数中调用 `WaitOneFrame()`。使用定长栈微线程时，将无法使用大量的栈，即使是临时使用。

这是一个很大的优点。如果 AI 例程被迫使用很浅的小型栈，则最终的 AI 将头脑简单，不能深入地思考问题。

### 3.2.4 并发症

#### 1. 装载和保存游戏

有时候，编译器将为每个函数生成栈帧，以便更容易地对局部变量和参数进行寻址。在栈上，这些栈帧以链表的形式连接到一起。换句话说，典型的栈包含指向自己的指针。因此，不能在其他地方使用栈映象（stack image）。

栈还包含返回地址——指向代码的指针。因此，将栈缓冲区保存到磁盘中后，如果重新编译了代码——所有的代码字节都将移动，则不能装载栈缓冲区，并指望它管用。

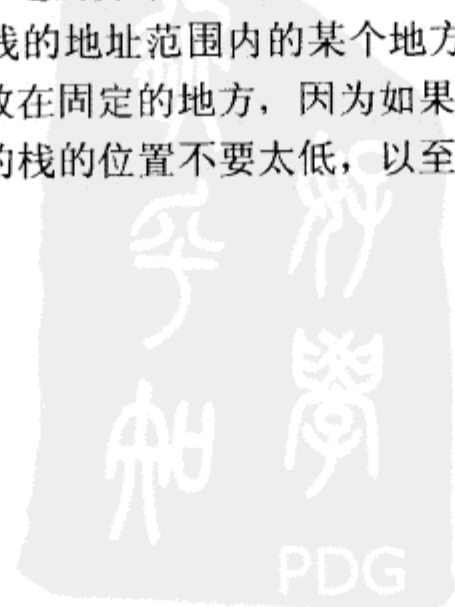
虽然对包含指向本身及代码的微线程栈的问题进行了处理，但栈将包含局部变量，而其中的一些局部变量为指针。当结构中包含指针时，只要小心谨慎，仍可以保存并恢复这种结构，但对于微线程栈，您不知道指针在哪里。小心地使用 C++ 指针能够处理这种问题，但非常复杂。

因此对于使用微线程栈的游戏，装载和保存是一个问题。这使得它只能用于不需要装载和保存的游戏以及控制台机器中。在控制台机器中，可以恢复保存时的内存布局。

#### 2. 结构化异常处理（structured exception handling）

Win32 结构化异常处理是调试策略的重要组成部分，因为它确保详细的崩溃信息能够被记录下来。它还被用于实现 C++ 异常处理。然而，如果不小心，对于发生在微线程中的异常，其结构化异常处理程序将不会被调用。这是因为发现不可信的指针（如位于当前 OS 线程已知的栈范围外的栈指针）后，遍历异常处理程序链表的 OS 处理程序将放弃[Pietrek97]。

如果将临时栈放在实际栈的地址范围内的某个地方（实际使用的地址下面），可以避免这种问题。记住，必须将它放在固定的地方，因为如果其位置发生移动，栈上的链表将无法工作，因此必须确保主线程的栈的位置不要太低，以至于被微线程栈所覆盖。



### 3. OutputDebugString



在 Windows NT 上, 如果调用 `OutputDebugString()`, 且不在调试器中运行时, 程序可能由于前面指出的结构化异常处理问题而退出。这很容易处理, 只需将栈放在合适的位置或使用附带光盘中的 `OutputDebugStringW95()` 函数即可。该函数检测是否是在调试器中运行, 且仅在安全时才使用 `OutputDebugString()`。它还检查 DBWin32——一个免费的调试输出监控程序, 并在可能的情况下同其交流。

#### 3.2.5 结论

为很多独立实体编写 AI 代码时, 线程是一种比较简单的方法。微线程让我们能够使用线程实现游戏实体, 而不会像其他线程方法那样带来较高的内存和 CPU 成本。在任何 CPU 体系结构中, 微线程实现起来都很容易, 只需很少的汇编语言代码[Keppel01]。



附带光盘中有一个用于 Win32 的微线程实现范例, 另外还有一个使用微线程的简单游戏, 以及一个将微线程同 Fiber 和 OS 线程进行比较的测试应用程序。

人们还使用多种用于游戏的脚本语言实现了微线程, 如 Lua[Lua01]、SCUMM (一种专用的 Lucas Arts 引擎, 也被 Humongous Entertainment 所采用) 和 Python[Tismer01]。

微线程“看门人”能够清理您的游戏代码。

#### 3.2.6 参考文献

[Brooks95] Brooks, Frederick P., Jr. *The Mythical Man-Month : Essays on Software Engineering, Anniversary Edition*, Addison-Wesley, 1995.

[keppel01] keppel.David, “QuickThreads distribution,” available online at <http://www.mit.edu/afs/sipb/project/scheme/src/guide-1.3/qt/>, June 5, 2001.

[Lua01] “The Programming Language Lua,” available online at [www.tecgraf.puc-rio.br/lua/about.html](http://www.tecgraf.puc-rio.br/lua/about.html), February 22, 2001.

[Pietrek97] Pietrek, Matt, “A Crash Course on the Depths of Win32 Structured Exception Handling,” *Microsoft Systems Journal* (Jan 1997).

[Tismer01] Tismer, Christian, “Stackless Python,” available online at [www.stackless.com/](http://www.stackless.com/), February 23, 2001.



## 3.3 使用微线程管理 AI

---

Simon Carter, Big Blue Box Studios

scarter@bigbluebox.com

正如前一篇文章中讨论的，游戏中的 AI 通常是通过某种状态机实现的。对于 AI 来说，状态机体系结构有很多优点。其中最重要的是，系统能够在任何特定的状态挂起，对于任何希望有多个 AI 实体的游戏来说，这都是至关重要的。然而，传统的状态机实现常常混乱、不直观、容易出现 bug、难以调试、难以理解。微线程提供了一种更精致的实现状态机的方式，并使得 AI 系统的健壮性和扩展性都非常高。本文介绍如何实现这样的系统，并充分利用其灵活性。

### 3.3.1 拼凑

---

微线程让我们能够使用常规的、日常的编程惯例来编写状态机。大部分“戏法”是在后台进行的，这让我们能够尽可能精致地设计 AI 系统，无需因后台体系结构方面的问题而分心。虽然很自由，但当这么多的限制烟消云散时，却可能难以知道从何处着手。

游戏中优秀的人工智能只是将构成复杂 AI 系统的内容组织成易于管理的块，并尽可能提高这些“模块”的直观性和可重用性。一个重要的设计决策是这种模块化的粒度。例如，传统的状态机常常将每种状态作为一个模块，这种非常低级的粒度使得系统笨拙而难以使用。微线程让我们能够自行选择模块的粒度，选择从概念上说是最合理的粒度水平。

从设计的角度说，将复杂的 AI 分解为“行为”单元要好得多。在这里，“行为”指的是一整套检测和行动，用于确定实体对特定刺激的响应，例如，饥饿、战斗、恐惧等。实体的行为越多，AI 将越丰满、变化越多端。理想的设计方案是，能够将不同的行为组同不同的实体关联起来。另外，还必须能够重用不同类型的实体共同的行为，从而最大限度地减少重复的代码量。

基于行为单元的系统让我们能够使用这些不同的模块构建出灵活的 AI，从而使用可重用的行为快速创建出不同的实体“大脑”。

### 3.3.2 良好的行为

---

下面是一个行为模块的接口：

```

class CAIBehavior
{
    CAIEntity* PEntity;
public:
    CAIBehavior(CAIEntity* pentity);
    virtual bool IsRunnable(void) = 0;
    virtual void Update(void) = 0;
    virtual void OnActivate(void) = 0;
    virtual void Cleanup(void) = 0;
};

```

**IsRunnable** 方法负责指出是否具备采取该行为的条件，对于“饥饿”行为，条件为附近是否有食品和实体是否感到饥饿。每当特定的行为激活后，**OnActivate** 都将被调用；而每当行为不再处于活动状态时，**Cleanup** 都将被调用，该方法负责确保实体被恢复到该行为激活前的状态。

不要对此过于担心，因为仅在后面组织多种行为时，才需要这些方法。就现在而言，对于单个行为单元，最重要的方法是 **Update**，该方法完成 AI 中的具体工作。

通常，AI 行为的特定模块只不过是一个条件测试树，将被分解为不同的行动（action）。例如，想想饥饿的动物可能做些什么。

```

CFood* pfood = FindPFood(PEntity->GetPos());
if(pfood!=NULL){
    //Move toward the food until we are near it
    while(!PositionsAdjacent(PEntity->GetPos(), pfood->GetPos()))
        PEntity->MoveOneStepTowardsPos(pfood->GetPos());
    PEntity->EatFood(pfood); //Eat the food.
}

```

这里，饥饿的动物将寻找食物，如果发现食物，则每次向它移动一步，直到靠近为止。靠近食物后，则食之。不幸的是，如果我们运行上述代码，游戏将在 **while** 循环中停顿不前，这正是微线程力挽狂澜的时机。假设该行为类在微线程中运行，则只需添加下述神奇的代码行：

```

while(!PositionsAdjacent(PEntity->GetPos(), pfood->GetPos())){
    PEntity->MoveOneStepTowardsPos(pfood->GetPos());
    MicroThreadSleep(); // a call to suspend the micro thread.
}
PEntity->EatFood(pfood);

```

转眼之间，简单的 AI 处理便变成了一个状态机！现在，每走一步后，AI 都将其处理挂起，将控制权返回给主游戏。

### 3.3.3 了然于胸

虽然看起来我们并没有取得多大成就，但我们现在能够创建单个 AI 行为模块了。然而，为使实体的生活丰富多彩、变化多端，我们需要对多个相互冲突的行为进行组织，并设置它

们的优先级，以便每次只有一种行为处于活动状态。这需要大脑类（**brain class**），如下所示：

```
class CAIBrain
{
    MicroThreadInfo*      PMicroThread;
    std::list<CAIBehavior> Behaviors;
    std::list<int>         BehaviorPriorities;
    CAIBehavior*          PActiveBehavior;
    int                   ActiveBehaviorPriority;
public:
    void AddBehavior(int priority, CAIBehavior& behavior);
    void Update(void);
};
```

游戏中的每个 AI 实体都将有自己的大脑，而每个大脑就是一个微线程。使用 **AddBehavior** 方法，不同类型的实体可以将不同的行为加入到行为集中。另外，每种行为都被指定一个优先级，用于帮助确定应采取哪种行为。大脑类的 **Update** 方法负责控制微线程，后者负责执行活动行为的 **Update** 方法。

然而，在此之前，我们必须确保有一个处于活动状态的行为。当实体无所事事时，我们需要遍历所有可能的行为，并从中选择一个通过 **IsRunnable** 测试的优先级最高的行为。如果该行为有专门的初始化代码，则对其调用 **OnActivate** 方法，并将其设置为活动的行为。该行为结束后，我们调用 **Cleanup** 方法，以运行其析构（**uninitialization**）例程。然后重复上述过程。

这样，大脑将确保总是有合适的行为在进行。实际上，给所有实体都指定某种总能通过 **IsRunnable** 测试的空转行为通常是个不错的主意，这样实体看起来好像总是在做事。

这样做有一个小小的并发症。假设因为万籁俱静，实体决定执行其 **Sleep** 行为几分钟，虽然其 **DefendSelf** 的优先级更高，但在其 **IsRunnable** 方法中没有发现进攻者。不幸的是，当它睡眠时，该实体遭到了不道德的敌人的攻击。使用前面介绍的系统，该 AI 实体将在睡眠中被残忍地杀害，因为它仅在当前行为结束后才会选择另一个活动的行为。

我们需要做的是，定期地查看实体行为清单，看有没有通过 **IsRunnable** 测试、且优先级高于活动行为的行为。如果找到这样的行为，则结束当前行为——这里为 **Sleep** 行为，并进入新的高优先级的行为——**DefendSelf** 行为。在这个例子中，就是要在进入新的、更为暴力的行为之前，调用 **Sleep** 行为的 **Cleanup** 方法，以确保实体被唤醒。这个系统的基本部分的代码如下：

```
void Update()
{
    CAIBehavior* pending = NULL;
    if(PActiveBehavior) {
        pending = FindNextBehavior(ActivePriority);
    }
    if(!pending && !PActiveBehavior)
        pending = FindNextBehavior(-1);
    if(pending) {
        if(PActiveBehavior)
```

```

        TerminateActiveBehavior();
        PActiveBehavior = pending;
        PActiveBehavior->OnActivate();
    }
    if(PActiveBehavior)
        SwitchToMicroThread(PMicroThread);
}

void FindNextBehavior(int priority)
{
    //Find a higher priority behavior that passes IsRunnable
}

static void MicroThreadFunction(void* pcontext)
{
    CAIBrain* pthis = (CAIBrain*)(pcontext);
    while(!pthis->TerminateThread){
        if(pthis->PActiveBehavior){
            pthis->ActiveBehaviorRunning = true;
            pthis->PActiveBehavior->Update();
            pthis->ActiveBehaviorRunning = false;
        }
        MicroThreadSleep();
    }
}

```

### 3.3.4 并发症

#### 1. 消亡后的通告

至此，一切都好像很简单。大脑运行微线程，而后者调用行为的 **Update** 方法。然而，跨越多回游戏进行处理时，存在很多并发症，我们必须对其进行处理。回到前面关于饥饿的代码，其中隐去了多个难点。首先是跟踪实体消亡的问题：由于代码现在跨越多回游戏，因此当我们达到我们跟踪的食物所在的地方时，该食物可能已被他人吃掉了。如果食物被吃掉后，只是将其从游戏中删除，则代码将崩溃，因为指针将指向无效的数据。

根据我们如何在游戏中处理实体的消亡，对这种问题的处理方式有很多种。从根本上说，我们必须确保在多回游戏之间，指向游戏实体的句柄是合法的。另外，我们还必须能够查询我们跟踪的实体是否已经消亡。一种解决方案是，让永久性句柄为标识符数字而非指针，且对于游戏中创建的每个实体都是惟一的。这样，当实体消亡时，游戏代码只需指出句柄无效即可。不幸的是，每当需要访问实体的任何部分时，都必须将句柄转换为指向实际数据的引用，会导致代码庞大、低效。







一种更好的方法是，编写一个特殊的智能指针系统，由它自动处理这些问题，Scott Meyers[Meyers96]提供了一整套编写智能指针的指南。从概念上说，传统的智能指针通过某种形式的引用计数，对其指向的对象有一定程度的所有权。我们的要求则相反——指针将自身的所有权交给它指向的对象。只要有跟踪指针指向某个对象时，便向对象注册自己。这样，当对象消亡时，便能够向指向它的所有指针通告，并将这些指针设置为空。采用这种方式后，在 AI 代码中只需在使用指针前检查它是否为空即可。有关如何实现跟踪指针系统的范例代码，请参阅附带光盘。

## 2. 消亡后的清理工作

还必须处理另一个问题。如果我们运行的 AI 所属的对象消亡或被中断，将发生什么情况？在这种情况下，需要停止活动的行为，删除在堆上创建的所有对象，并尽快地退出 AI 代码。同样，根据您的偏好，处理这种问题的方法有多种。

一种妥善的方法是使用 C++ 异常处理功能。引发异常后，幕后的 C++ 神奇之手将对栈进行清理，并沿异常处理层次结构向上传递控制权，这正是我们所需的功能。大脑类所需做的只是捕获异常，并调用行为的 Cleanup 方法。然而，这是一种极具重量级的方法，有些微线程实现对异常的处理并不是太友善。

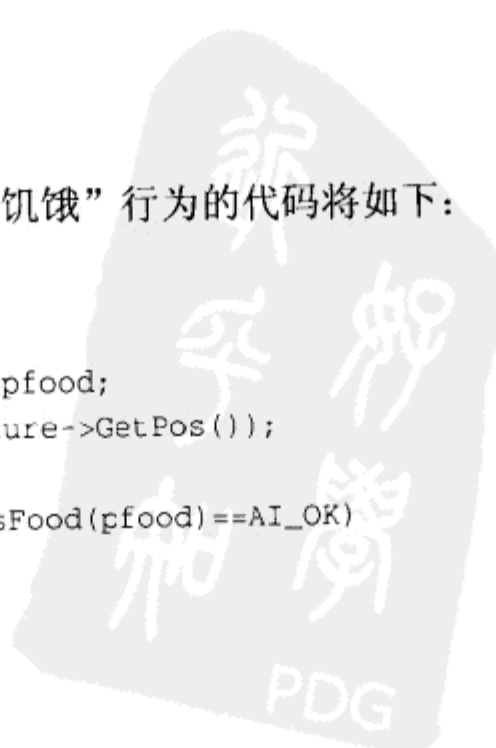
一种更简单（也更具侵略性）的方法是，定期地检查是否被终止，如果是则返回。这确实要求要仔细地组织 AI 代码，但这也是一种非常轻量级的解决方案。稍微动动脑筋，甚至可对其进行封装，使其侵略性降低到最小。在大多数情况下，为确保安全，大脑类将不断轮询微线程，直到获悉行为确实已停止运行为止，如果您以前使用过线程，则应熟悉这种技术。

```
void CAIBrain::TerminateActiveBehavior()
{
    if(PActiveBehavior){
        PActiveBehavior->SetTerminateFlag();
        while(ActiveBehaviorRunning)
            SwitchToMicroThread(PMicroThread);
        PActiveBehavior->Cleanup();
        PActiveBehavior->ClearTerminateFlag();
    }
}
```

## 3. 动作发出通告

进行上述修改后，前面关于“饥饿”行为的代码将如下：

```
void Update()
{
    CTrackingPointer<CFood> pfood;
    pfood = FindPFood(PCreature->GetPos());
    if(pfood!=NULL){
        if(ActionMoveTowardsFood(pfood)==AI_OK)
```





```

        PEntity->EatFood(pfood);
    }
}

EAIReturn ActionMoveTowardsFood(CSmartPointer<CFood> pfood)
{
    while(pfood!=NULL && !PositionsAdjacent(PEntity, pfood->GetPos())){
        if(Terminated())
            return(AI_TERMINATE);

        PEntity->MoveOneStepCloserToPos(pfood->GetPos());
        MicroThreadSleep();
    }
    if(pfood!=NULL)
        return(AI_OK);

    return(AI_FAIL);
}

```

这里使用一个特殊的跟踪指针来跟踪食物，以应付其消亡的情况。ActionMoveTowardFood方法每向食物前进一步后，便要求大脑将微线程挂起。如果动作（action）发现自己被行为的大脑终止，则返回一个值，让调用它的 AI 行为知道它应干净利索地退出。同样，如果食物消亡，它将返回一个编码，将这一点告知行为。另外，对跨越多回游戏存在问题的过程的处理代码被封装到一个独立的方法中。以这样的方法组织代码有很多重要的优点：

- AI 实际上是一棵条件测试树，它最终被分解为实际要做的事情。这些“动作”只是跨越多回游戏的代码片段，因此将其同系统的其他部分隔离开来是合理的；
- 由于只有动作需要将处理挂起，因此只有它们需要检测终止情况。将挂起代码和终止代码放在一起将更为整洁，并降低了顾此失彼的可能性；
- 动作函数可放在行为基类中，并在不同的行为间重用。

#### 4. 扩展

这里描述的系统较松散，这是作者有意为之的，因为作者只想提供一种这样的思想，即微线程能够提供给 AI 精致的体系结构。可以在其中添加任意的 AI 技巧，包括：

- 给每种行为指定一个字符串名称，以便调试时显示。
- 使用外部数据脚本语言组装大脑，让非程序员能够创建和具体化实体。有关这个主题，请参阅本书中 Scott Bilas 撰写的文章。
- 使用《游戏编程精粹 1》中描述的消息系统，让 AI 实体能够查询和发送有关有趣事件的信息[Rabin00]。

### 3.3.5 结论

微线程为编写 AI 提供了巨大的灵活性。正确地利用这种灵活性，可以快速而轻松地创建出合理、快速、轻量级、没有 bug 的 AI。与其他方法相比，这种系统的优点包括：

- 以直观的方式将行为划分为模块。通过将其加入到不同的大脑中，可在不同类型的实体间重用它们。
- 不会像其他状态机实现中常见的那样，沿状态之间混乱的链接随意跳转。这样，调试时您能够明白整个条件树，确定导致问题的原因，而无需跟踪不同的状态。
- 程序员可以使用其喜欢的技术来编写代码，而不受制于严格的体系结构。
- 在游戏进入到下一回时，可将实体特有的数据及其行为存储在行为类中。

### 3.3.6 参考文献

---

[Meyers96] Meyers, Scott, "Smart Pointers," *More Effective C++*, Addison Wesley, 1996.

[Rabin00] Rabin, Steve, "Designing a General Robust AI Engine," *Game Programming Gems*, Charles River Media, 2000.



## 3.4 一种 RTS 命令排队体系结构

---

Steve Rabin, 任天堂(美国)公司

steve\_rabin@hotmail.com

**实**时策略游戏有一种独特的用户交互方法。通过使用鼠标, 玩家能够给个体或个体组下达多个命令。这种交互已成熟多年, 每种新 RTS 都是在以前设计的基础上构建的, 并对其进行了改进。一种进展最快的设计是命令排队技术。本文介绍这种交互方法的工作原理以及如何直接将它集成到底层的 AI 体系结构中。

### 3.4.1 RTS 命令

---

RTS 游戏的基本用户界面涉及到选择个体, 并命令他们执行某项任务, 如向敌人发起攻击或移动到特定的地方。这是一种非常简单的概念, 不需要任何游戏的 AI 体系结构就能实现。然而, 讨论更复杂的组合之前, 必须考虑大多数 RTS 游戏中都有的简单命令。下面是常见的 RTS 命令列表:

- 向建筑物或敌人发起攻击;
- 构筑建筑物或制造武器;
- 移到某个地方;
- 巡逻(patrolling)到某个地方;
- 收集资源(食物、原材料、能源);
- 研究某种技能;
- 修理设备或建筑物;
- 回收设备(回收报废设备中的材料/能源);
- 保卫设备或建筑物(向来犯之敌发起攻击);
- 占据某个位置(向一定范围内的敌人发起攻击, 但不移动);
- 停止;
- 自毁。

### 3.4.2 命令排队

---

玩 RTS 游戏时, 您将大量的时间用于命令个体移到什么地方。不幸的是, 在大多数游戏中, 寻径都不完美, 玩家通过设置中继点(waypoint)来规划路径将会有所帮助。中继点是个体中排队等候的连续移动命令。玩家在单击地面的同时按下特定的按钮(如 Shift 键), 将每个中继点的移动

命令排队。

中继点排队是重要的一步，它打开了通向强大界面系统的大门。如果能够对中继点排队，为何不让玩家对任何命令组合进行排队呢？实际上，您可以命令个体向敌人发起攻击、修补城墙并建立炮塔，而无需等待这些任务完成。另外，玩家稍后可以将更多的命令加入到队列末尾。通常，这种概念被统称为排队。

这里的诀窍在于，将每个命令视为一项任务，并将个体的大脑视为一个任务队列。个体总是处理队列开头的任务。任务完成后将被删除，然后下一项任务执行。当没有任务需要处理时，个体应有一个默认的空闲任务。图 3.4.1 是一个个体大脑队列的例子。

图 3.4.1 是将命令攻击、移动、修理等命令进行排队后的结果。其中每个命令都有与之相关联的数据，如攻击目标或移动到什么地方。新任务被加入到队尾，但位于默认任务之前。默认任务总是位于最后，它永远不会结束，也不会被删除。

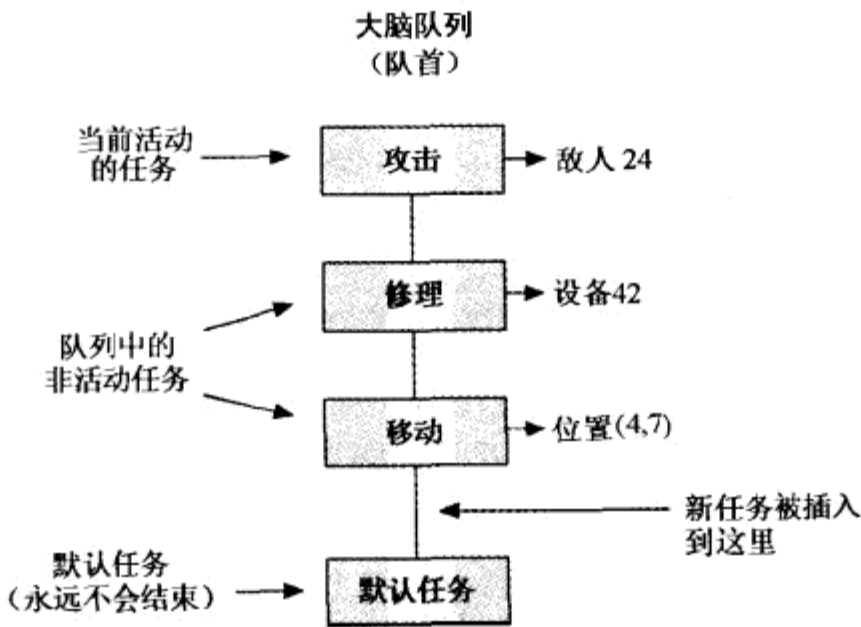


图 3.4.1 一个大脑队列形式的 AI 任务列表

如果玩家命令个体执行一项新任务，而没有按下“排队”按钮，则所有的任务都将被删除，同时新任务将被加入到队列中。因此，很容易用单个新任务替换队列中的任务。

在这种体系结构下，一种命令行为是，允许玩家查看命令队列，方法是选中个体的同时按下“排队”按钮。屏幕上的结果是一些箭头，指出个体将去哪里以及打算干什么（如建造建筑物）。对此，一种精致的实现方式是，让每项任务将其作用绘制到屏幕上。这样，玩家便能迅速明白个体队列中的任务以及还应添加哪些任务。

这种简单的体系结构并不知道是谁将命令加入到队列中的。显然，玩家可以将命令加入到队列中，但游戏也可以预先为非玩家角色（Non-Player Character, NPC）装载命令，这些命令可能是巡逻到某个地方或保卫某个地方。AI 还能够就攻击哪些个体作出高级决策，为此，只需将相应的命令加入到队列中即可。甚至可以允许通过网络传递的分组将命令加入到个体的队列中。这是一个很精致的系统，有很大的灵活性。

3.4.3 循环命令

巡逻（patrolling）是一个循环命令，它给前面介绍的系统带来了一些有趣的结果。当玩家命令个体巡逻到某个地方时，该个体将记住其原来的位置，并不断地在这两点之间来回走动（直到发现敌人为止）。玩家也可以将多个巡逻中继点加入到队列中，这样个体将不断地在这些中继点之间循环走动。图 3.4.2 是一个 3 点的巡逻，这是通过单击两次鼠标启动的。

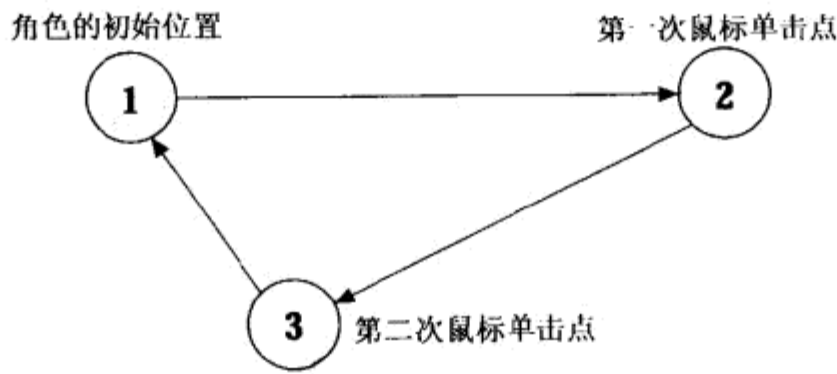


图 3.4.2 个体的巡逻路径

1. 将第一个巡逻点加入到队列中

玩家单击第一个巡逻点实际上将把两个巡逻命令加入到队列中。这是因为玩家的意图是，要求角色移动鼠标单击的位置，然后返回其原来的位置，并不断重复这一过程。图 3.4.3 是玩家发出一个巡逻命令后的大脑队列。

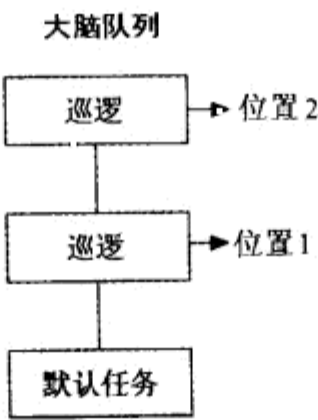


图 3.4.3 玩家发出第一个巡逻命令后的大脑队列

有趣的是，除了循环到队尾外，巡逻命令同移动命令完全相同。因此，只需在移动任务中设置一个“循环”标记，使之成为巡逻任务，我们的排队系统便能很好地工作。巡逻任务结束后，将被放到队尾，而不像移动任务那样被删除。图 3.4.4 说明了图 3.4.2 所示的 3 点巡逻的几次迭代。

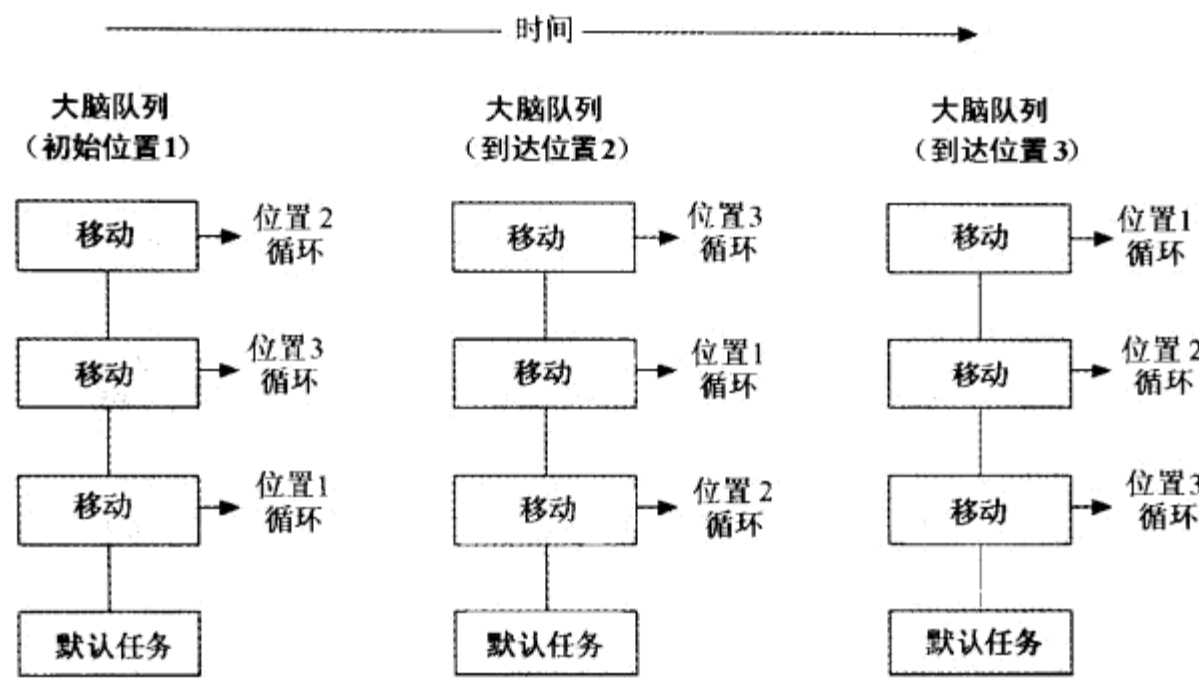


图 3.4.4 巡逻路径的三次迭代

2. 将其他命令加入到队列中

队列中包含巡逻任务后，如果玩家要将其他命令加入到队列中，将引发一些微妙的问题。虽然下面的内容有些主观，但试验表明，这与大多数玩家所期望的很接近。

第一个问题是将其他巡逻命令加入到队列的什么位置。玩家可能希望将它们加入到第一组巡逻点的后面（而不管个体当前位于大脑队列的什么位置）。这很重要，因为将第二组巡逻点加入到队列中，个体可能位于巡逻路径的任何位置。

解决方案是，对最后一个被加入到队列中的巡逻点进行标记，这样便可以将新的巡逻命令加入到它后面。加入新的巡逻命令后，移动标记“last”。图 3.4.5 是一个 3 个巡逻命令依次被排队的例子。

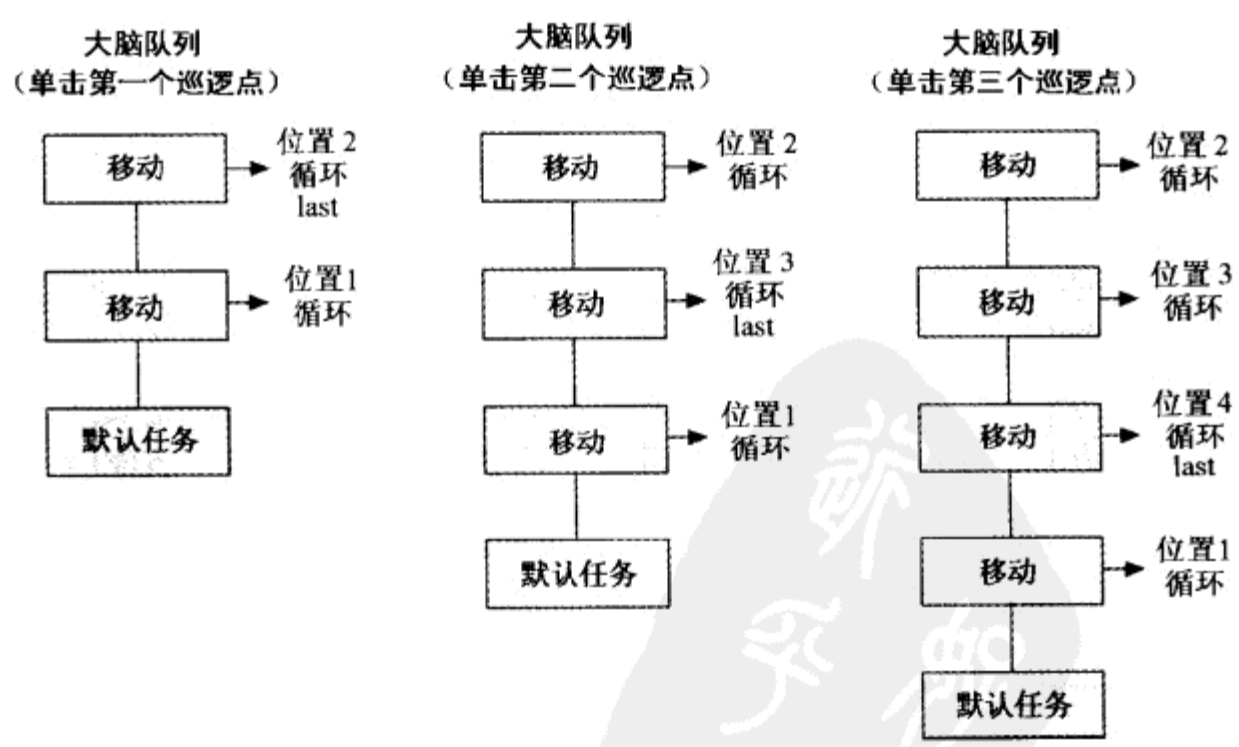


图 3.4.5 依次对 3 个巡逻命令进行排队



对巡逻 (Patrolling) 命令进行排队时涉及到的第二个问题是, 将非巡逻命令加入到队列中。通常, 玩家希望完成当前巡逻点后, 立刻执行新加入的命令。这有些棘手, 因为新命令必须加入到当前巡逻命令 (如果有的话) 后面, 且位于其他非循环命令之后。图 3.4.6 是一种对多个新的非巡逻命令进行排队的情况。

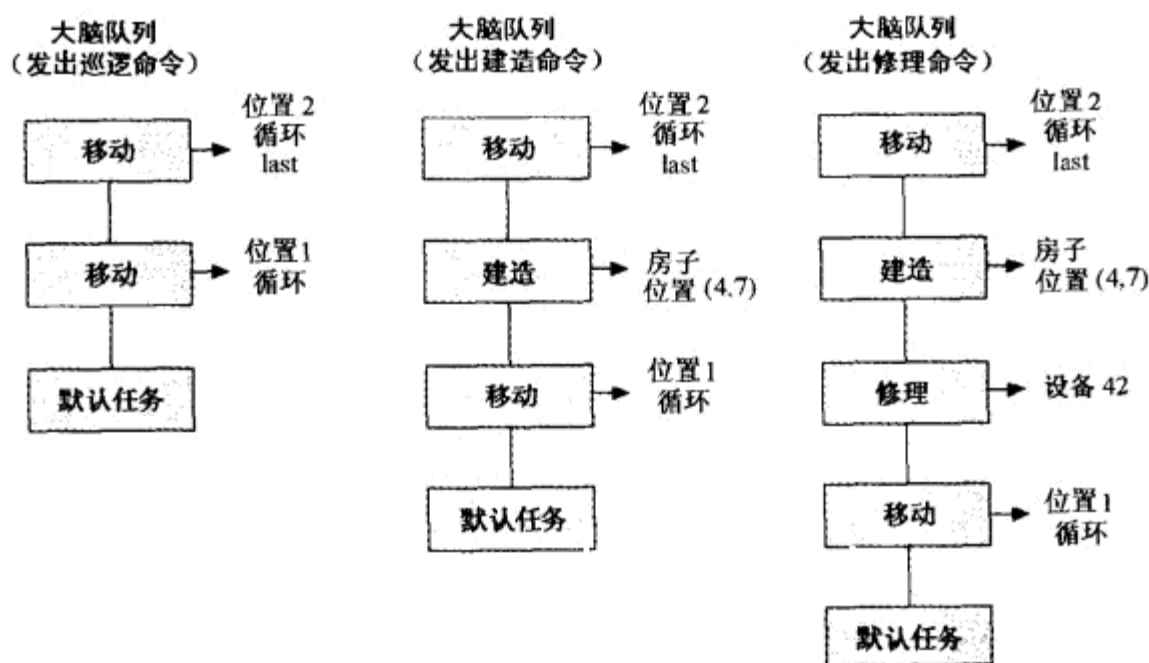


图 3.4.6 将其他命令加入到包含巡逻任务的队列中

从图中可知, 巡逻命令给命令排队概念带来了多种挑战, 其关键在于要按用户期望发生的情况来实现。用户通常对系统的工作方式有一定的模型, 了解这种模型并非一项容易的任务, 但这样将让您知道如何设计命令排队系统的交互和行为。

然而, 了解用户的思维模型是双向的。给玩家提供及时的反馈, 让他们知道其输入的结果以及输入将如何被解释同样重要。这样, 如果用户的模型不正确或存在缺陷, 则可以迅速调整, 以获得正确的模型。讨论思维模型的精彩图书之一是 Donald Norman 编写的 *The Design of Everyday Things* [Norman90]。

### 3.4.4 结论

命令排队现在是一种标准特性, 没有它, 任何 RTS 游戏都将寸步难行。通过使用大脑队列体系结构来存储任务, 消除了命令排队系统的很大一部分复杂性。另外, 可以将大脑队列视为一个简单的任务列表, 也可以将每项任务分别封装为一个知道如何完成其工作的 AI 系统。无论采用哪种方式, 玩家都能将其要发出的命令进行排队, 并轻松地使用这个简单而强大的界面管理数以百计的个体。

### 3.4.5 参考文献

[Norman90] Norman, Donald, A., *The Design of Everyday Things*, Currency/Doubleday, 1990.

## 3.5 一种基于分片的高性能视域和搜索系统

Matt Pritchard, Ensemble Studios  
mpritchard@ensemblestudios.com

在策略游戏领域，在目标识别和探测时，常常遇到视域（Line-of-Sight, LOS）和战争烟雾（Fog-of-War, FOW）的概念。在传统的基于回合（turn-based）的策略游戏中，采用蛮干方法重建玩家已探明的视域地图就足够了。然而，将这种方法用于实时策略（RTS）游戏时，尤其当游戏世界的规模和实体数目增大时，很快便显露出其性能方面的问题。很多商业 RTS 游戏（其中一些还非常成功）仍然使用这种完全重建的方法。不幸的是，它们不得不在性能上进行折衷，如放弃战争烟雾或不是每回都对其进行更新，允许显示的地图不那么精确。本文介绍一种高效的、用于基于分片的游戏中的玩家可见性（visibility）系统，它在提供快速搜索以及其他游戏特性的同时，最大限度地降低了性能影响。

### 3.5.1 概述

第一个假设是，在程序的内部，游戏世界的各个部分是用分片（tile）表示的，分片对应于一个二维阵列。

· 大多数实时策略游戏都使用矩形地图，但很容易对该实现进行修改，以适应基于六边形的世界布局。

我们的玩家可见性系统的目标如下：

- 在任何时候，玩家探明的、可见的、不可见的分片都必须是完全精确的；
- 个体能够基于 FOW 和可见性，快速地搜索其他个体或物体；
- 系统必须同时支持多达 16 位玩家，并允许玩家之间共享任何信息；
- 系统必须有较好的伸缩性，能够支持更多的个体、更大的地图和更大的搜索半径。

### 3.5.2 定义

- 分片（Tile）：游戏世界中最小的、离散的方形或六边形部分；
- 世界（World）：游戏发生的整个区域；在内部为一个二维分片阵列，也称为地图；
- 玩家（Player）：控制所有个体的人或计算机实体；每个玩家的可

见度信息都不同;

- 个体 (unit): 玩家控制的任何游戏实体, 可以是不能移动的;
- 视域 (Line-of-Sight, LOS): 个体可见的区域;
- 视域半径: 个体的 LOS 范围, 单位为分片;
- 未探明的分片 (Unexplored Tile): 在玩家个体的 LOS 内见不到的分片;
- 不可见的分片 (Fogged Tile): 已被探明但不在玩家任何个体的 LOS 内的分片;
- 可见的分片 (Visible Tile): 在玩家的一个或多个个体的 LOS 内的分片;
- 战争烟雾 (Fog Of War, FOW): 是这样一种概念, 即如果已探明的分片当前不在玩家的个体的 LOS 内, 则该玩家将看不到其他玩家的位于该分片内的个体。

在 RTS 中, 我们假设对于任何一位玩家, 开始时地图都是未被探明的。当分片进入个体的 LOS 内后, 这些分片对于个体所属的玩家而言将是已探明的。当已探明的分片不再位于该个体的 LOS 内时, 则对于相应玩家而言, 该分片将是不可见的, 但不会回复到未探明的。需要指出的是, 分片是未探明的、可见的还是不可见的随玩家而异, 这很重要。

### 3.5.3 组件 1: 各个玩家的可见性地图

---

对于玩家可见性系统, 需要实现的第一个组件是游戏中每个玩家的可见性地图。这是一个非常简单的结构: 一个二维 byte 数组, 它与分片布局具有一一对应的关系。每个数组元素都包含是一个数字, 指出玩家有多少个个体能够看到该分片 (也就是说, 该分片位于多少个个体的 LOS 内)。

更新可见性地图很容易。当个体刚被创建或移动到新的分片位置时, 对于位于该个体视域内的所有分片, 其可见性计数都加 1。当个体被删除时, 则对于原来位于该个体视域内的所有分片, 其可见性计数都减 1。如果分片对玩家而言是可见的, 则相应的可见性地图元素的值为非零。然而, 如果元素的值为 0, 则无法确定相应的分片是未被探明的还是不可见的。为解决这个问题, 我们指定 0 表示不可见, -1 (存储为 byte 类型时为 255) 表示未探明。

不幸的是, 当未探明的分片对应的值被加 1 时, 将变为 0, 这表示分片是不可见的, 而这是不正确的。然而, 我们可以检测这种特殊情况, 确保如果加 1 后为 0, 则再加 1。这提供了一个方便的地方, 可以加入对探明的分片进行一次性处理的代码, 如将其加入到微型地图中、记录其类型或搜索其中的资源。由于大多数游戏都不会将分片回复为未探明的, 因此这种特殊情况不会出现在递减代码中。有必要指出的是, 元素的大小 (这里为 byte) 限制了可以同时看到某个分片的个体的最大数目, 这里为 254 个。

### 3.5.4 组件 2: LOS 模板

---

在大多数策略游戏中, 个体的 LOS 被定义为以个体为中心的圆形区域, 半径为分片数。计算这种形状的最简单的方式是看分片离个体的距离是否小于个体的 LOS 半径, 很多游戏都是这样做的。然而, 从性能的角度看, 这种方法的效率是非常低的。由于每一回中, 这种操作需要占用大量的时间, 因此需要对这种函数进行优化。

一种优化方法是, 对于游戏中可能使用的每种 LOS 半径, 预先计算相应的 LOS 区域。

然后, 将这种形状信息存储在一个模板结构中, 并使用不同的模板表示不同的 LOS 半径。

作者发现, 最佳的实现是将 LOS 区域存储为一系列的水平带 (strip), 这些水平带用起始位置、结束位置以及相对于个体位置的垂直距离表示, 并按从上到下的顺序排列。假设在内存中, 数组水平轴中的元素是线性存储的, 则水平地处理模板。对于位于边缘的个体, 则根据游戏地图对 LOS 模板的形状进行剪切, 这只需在外部循环中对起始和终止值进行定位 (clamping) 即可。下面是一个将个体的 LOS 模板加入到可视性地图中的函数的代码。

```
// This routine "explores" the game map in an area around the specified
// unit position using a line of sight template
// The template span data is an array of structs with 3 offset elements:
// {vertical position, horizontal start, horizontal end}

void VisibilityMap::AddLOSTemplate(int XPosition, int YPosition, LOSTemplate *template)
{
    int n, x, y, xStart, xEnd;

    for (n = 0; n < template->number_of_vertical_lines; n++)
    {
        y = Yposition + template->SpanData[n].Yoffset;
        if (y >= map_y_min_coord && y <= map_y_max_coord)
        {
            xStart = max(XPosition + template->lines[n].XStartOffset, map_x_min_coord);
            xEnd = min(XPosition + template->lines[n].XEndOffset, map_x_max_coord);

            for (x = xStart; x <= xEnd; x++)
            {
                if ((VisibleMap[y][x]++) == 0)
                {
                    ExploreTileForFisrtTimeHandler(x, y);
                    VisbleMap[y][x] = 1;
                }
            }
        }
    }
}
```

当玩家的个体从游戏世界消失时, 游戏将其 LOS 区域的可视性计数减 1。如果分片不位于玩家的任何一个个体的 LOS 内, 则计数为 0, 这表明玩家再也不能看到该分片, 也就是说该分片是不可见的。

当个体从一个分片移向另一个邻接的分片中 (这是一种非常常见的操作) 时, 将从原来的位置中删除其 LOS, 并将其加入到新的位置中。由于这两个函数调用通常位于个体的移动代码的前后, 因此可以进行另一种优化, 即将这两种操作合并为一个函数。新函数将以新旧位置为参数, 且只更新玩家的可视性地图中不同时递增和递减的部分。另一种情况是个体的 LOS 半径发生变化。在这种情况下, 将调用删除 LOS 的函数, 并以原来的半径为参数; 然后调用添加 LOS 的函数, 并以新半径为参数。只要正确地进行编写, 优化后的更新函数也能够处理这种情况。

使用 LOS 模板还有其他优点。首先, 可以为 LOS 半径相同但大小不同的对象创建不同的形状。小型游戏个体可能只占据一个分片, 而大型个体 (如固定的建筑物) 则可能占据多个相邻的分片, 其形状可能不是方形的, 而是矩形或椭圆形。位于占据单个分片的个体中央的 LOS 模板, 被用于大型个体时可能不再在中央。图 3.5.1 是用于两个大小不同的对象的 LOS 模板形状, 这两个 LOS 的半径都是 3 个分片。

使用模板的另一个优点是, 可以创建不对称的形状。图 3.5.2 是一个定向探照灯形状的例子。通过一整套旋转 LOS 模板, 可轻松地让探照灯个体进行全方位照射, 实现一种非常酷的游戏效果, 而几乎不需要进行专门的编程。

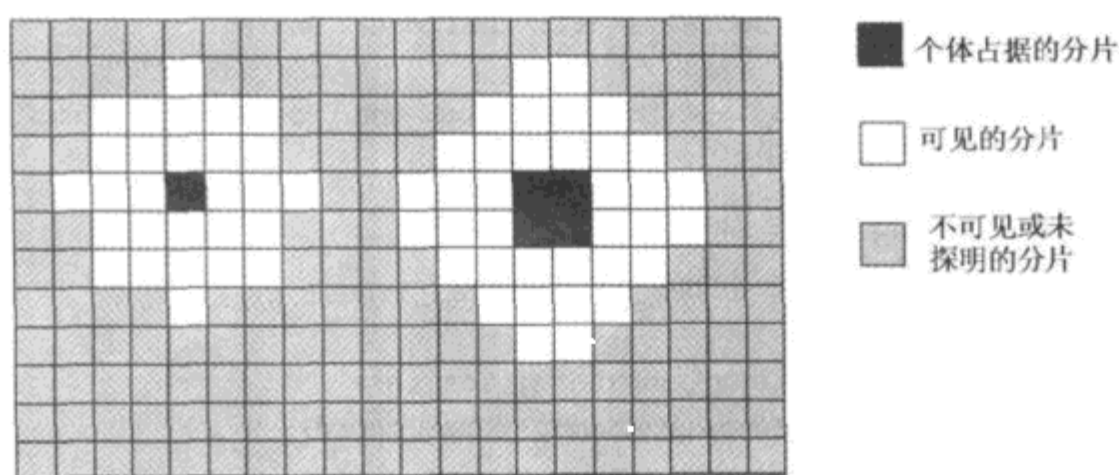


图 3.5.1 半径相同的对应于大小不同的个体的 LOS 形状

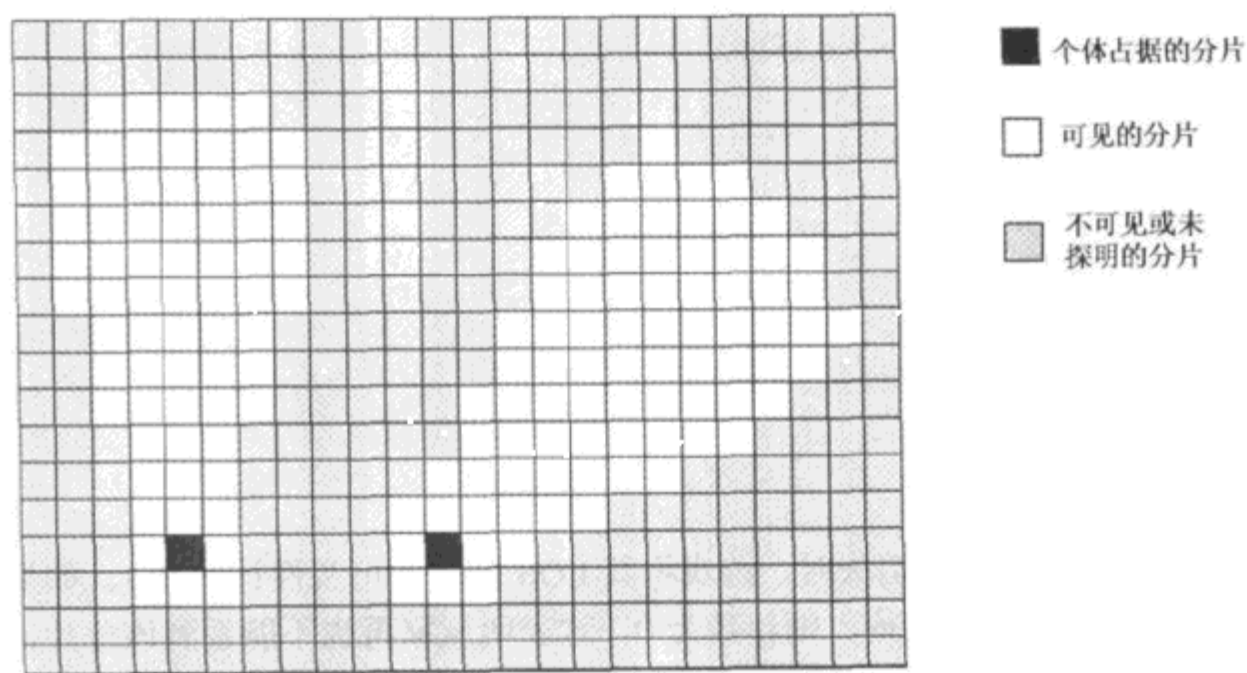


图 3.5.2 非圆形 LOS 区域, 有两种旋光的探照灯模式

### 3.5.5 组件 3: 合并的可视性地图

至此, 实现的效率提高了, 但这对其他目标并没有什么帮助。接下来的组件名为合并的可视性地图, 将把其他结构组合起来。该数据结构主要供其他游戏代码访问, 并给搜索函数带来很大的帮助。

与玩家地图一样, 合并的可视性地图也是一个二维数组, 其大小同分片网格相同。不同



之处在于，整个游戏中只有一个合并的可视性地图，其元素为 32 位的 DWORD，而不是 byte。鉴于其用途，将其声明为全局的，在整个程序中都是可用的是个不错的主意。

合并的可视性地图旨在在一个地方提供最新的、关于游戏中所有玩家的可视性信息。这只是通过对每个玩家使用一个 2 位的元素来实现的。其中一位指出分片对玩家而言是否已探明，另一位指出对玩家而言分片当前是否是可见的。这样，每个 DWORD 可提供 16 个玩家的数据。

在合并的可视性地图元素中，各个位的组织方式取决于用户如何实现。组织方式不会影响性能，因为所有的更新都是使用预先为每个玩家计算好的掩码值（mask value），对整个 32 位的元素执行二进制 OR、AND 或 XOR 运算的。

实际上，初始化组合的可视化地图时，所有的分片都被设置为未探明和不可见的。然后，当任何玩家发生下述事件时，对其进行更新：

- 分片首次被探明；
- 分片从未探明或不可见变为可见；
- 分片从可见变为不可见。

在显示游戏世界或其他结构（如雷达地图）时，每个玩家都有一个包含探明位和可视性位的掩码值。当多种函数遍历各个分片时，将组合的可视性图中该分片的值同当前玩家的可视性掩码进行 AND 运算。得到的结果指出了对指定的玩家而言，该分片的探明状态和可视性状态，然后代码根据结果进行相应的操作。

使用合并的可视性地图的第一个优点是，可以合并玩家的可视性掩码。这样可以实现各种游戏效果，如供队友共享视域和可视性以及监视其他玩家。随时可以添加或删除每个玩家的效果，只需更新玩家的可视性掩码并刷新屏幕即可。

### 3.5.6 改进搜索

直接的搜索方法需要查看 LOS 区域中每个分片中的个体。随着游戏中的个体越来越多，每回的搜索次数将增加。随着搜索半径的增加（如规定范围内的个体数或 LOS 增大），需要搜索的分片数将急剧增加。例如，对于单个攻击个体，如果搜索半径为 10 个分片，则搜索区域大约有 350 个分片。因此，这种直接方法的性能将随需要搜索的分片数线性降低。

当需要在个体的 LOS 内搜索敌人或其他感兴趣的东西时，使用合并的可视性地图的最大优点便显现出来了。与其对各个分片进行搜索，不如保存给定玩家的总体 LOS 中的其他玩家的个体列表。这正是合并可视化地图的用武之地。游戏中的每个个体通过访问合并的可视化地图来获悉其占据的分片。这样，每个个体就知道哪些玩家能够看到它。通过保存前一回中的这些信息，即使个体本身并没有移动，个体也将知道自己是否进入或离开了游戏中其他每个玩家的 LOS。另外，当发生变化时，个体可以将自己加入到其他玩家可见的个体列表中或从中将自己删除。除个体进入或离开另一个玩家的 LOS 外，更新开销只是对 DWORD 进行检查，每个个体每回一次。

可以根据对玩家而言个体是什么（如队友、战斗设备、基础设施），对个体将自己加入或从中将自己删除的列表进行分解。结果是，每个玩家都有一系列非常小的（甚至是空的）列表，其中包含指向其他玩家的个体的指针，这些个体当前位于该玩家的总体 LOS 中。

维护这样的列表后，便无需搜索大量的分片（其中大部分分片可能没有包含潜在的目

标)。相反,搜索变成只对一个小的、只包含潜在目标的列表进行扫描。这样搜索代码将更简单,同时可以更好的将列表放在高速缓存中。性能将提高多个数量级,尤其是当搜索涉及到很多远距离的个体时。

下面的代码说明了个体的更新如何处理 LOS 可视性中的变化。代码判断哪些玩家的可视性状态发生了变化,然后判断如何修改它们。

```
void GameUnit::TurnUpdateProcess(. . .)
{
    // Game specific unit processing code...
    . . .
    // Now we check to see if we've gone in or out of anyone's LOS
    DWORD CurrentlyVisibleTo = CombinedVisibilityMap[Yposition][Xposition];

    if ( CurrentlyVisibleTo != LastVisibleToValue)
    {
        // Get only the bits that have changed
        DWORD VisibilityChanges = CurrentlyVisibleTo ^ LastVisibleToValue;
        LastVisibleToValue = CurrentlyVisibleTo; // Save new value

        for (int playerNo = 0; playerNo < theGame->numOfPlayers; playerNo++)
        {
            DWORD PlayerMask = 0x0001 << playerNo; // bit mask for player

            // Check to see if our visibility for this player has changed
            if ((VisibilityChanges & PlayerMask) != 0)
            {
                if ((CurrentlyVisibleTo & PlayerMask) != 0)
                    AddUnitToPlayersVisibleList(playerNo, self);
                else
                    RemoveUnitFromPlayersVisibleList(playerNo, self);
            }
        }
    }
    // Continue with game processing
}
```

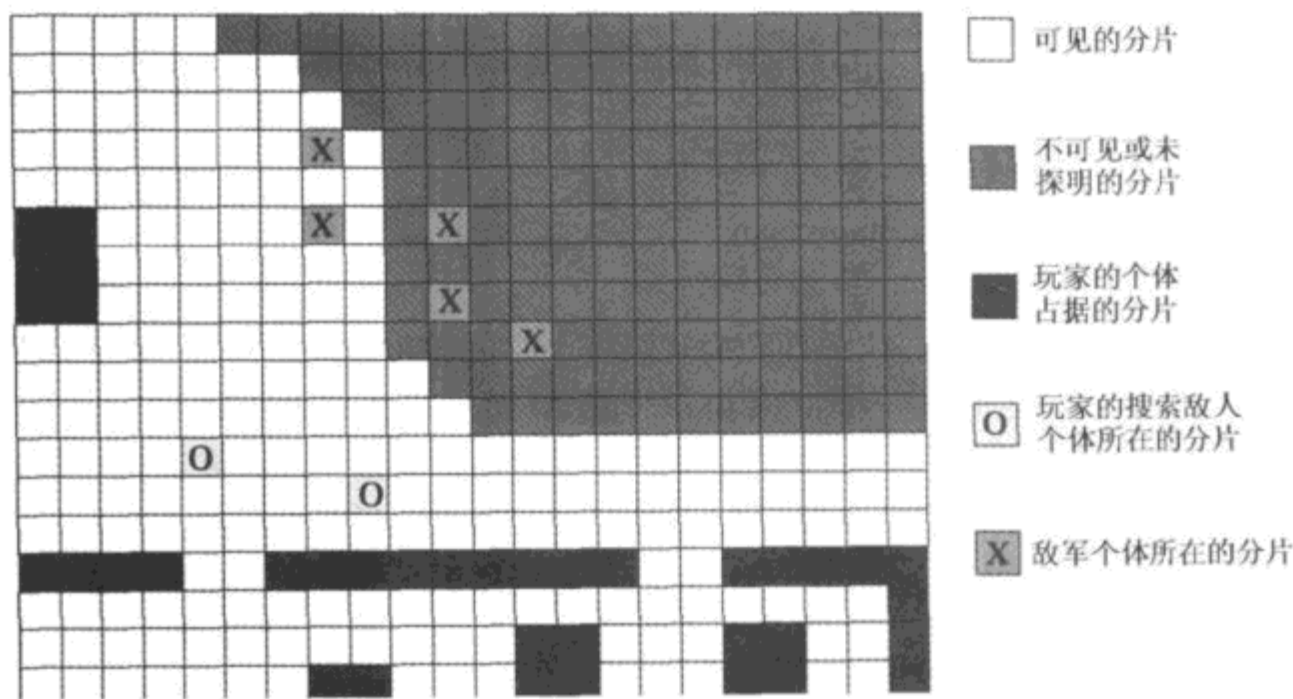
这种方法的另一个优点是,搜索将根据玩家的总体可视性进行,而不局限于个体的 LOS 半径。如图 3.5.3 所示,位于底部的玩家的每个个体对目标进行搜索(搜索半径大于其自身的 LOS 半径)并在左上方找出玩家看得见的敌军个体。然而,在右上方还有两个敌军个体,但它们却没有被发现,因为这些个体位于玩家看不到的分片上。根据合并的 LOS 并使用与 LOS 半径不相关的搜索半径,可以做出更为智能的、更人性化的决策。

同样的处理代码可用于占据游戏世界中多个分片的个体。惟一的区别是,当前的可视化值是通过合并的可视化地图中每个被占据的分片对应的值执行 OR 运算来得到的,而不是直接来自单个分片的值。这也使得揭示个体占据的整个地图区域更容易,即使玩家的 LOS 只覆盖了该个体的一个角。

合并的可视性地图使得易于实现的另一种游戏功能是蜃景(mirage)。蜃景是个体的幻影表示(ghost representation),它出现在另一个玩家不可见的区域中,表示该玩家在探明该区



域时看到的東西。能够生成场景的个体在离开另一个玩家的 LOS 时将这样做。另外，当场景（通过检查合并的可视性地图中的值）发现自己对另一个玩家而言完全可见时，将把自己删除。



3.5.7 结论

对事态重新进行思考总没有害处，即使它看起来简单明了。本文介绍的方法就是询问“为什么”，从而明白看起来不相关的系统可以相互借鉴。之所以能够在搜索功能方面获得重大改进，是因为有了合并的可视性地图。这之所以有效，仅仅是因为在整个游戏进程中，各个玩家的可视性地图总是最新的。因此，各个系统完成其特定的工作，而它们携手工作时，便能实现更多的功能并获得进一步的优化。



## 3.6 创建影响力地图

---

Paul Tozour

gehn29@yahoo.com

如果游戏 AI 主体在给定场境下总能做出合适的决策，则被认为是“智能的”。因此，可以说智能 AI 游戏的核心是决策。

日常经验告诉我们，要做出有效的决策，不仅仅要有最佳的数据，还需要以正确的方式表示这些数据。原始数据只有被转换为场境信息后才能发挥作用。以合适的方式表示数据可以揭示出底层相关的模式。

本文和接下来的一篇关于策略评估的文章将介绍这样一些技术，即让 AI 主体能够从战术和战略的角度深入了解游戏世界中的角色和当前的游戏状态，还将讨论与不同的游戏类型和虚拟环境相关的技术。

### 3.6.1 影响力地图

---

对于进行战术评估而言，建立影响力地图是一种宝贵的、得到实践检验的技术。影响力地图 (influence map) 最常被用于策略游戏中，但对于需要进行战术分析 (tactical analysis) 的很多其他游戏类型也很有用。建立影响力地图的概念是任何 AI 开发人员工具包不可或缺的组成部分。

影响力地图是 AI 主体关于游戏世界知识的空间表示，让计算机玩家 (computer player) 能够根据游戏环境的物理/地理表示获悉当前游戏状态的战术视图。影响力地图指出了计算机玩家的兵力部署位置、敌军的位置或最可能位于什么位置、玩家之间的“边界”、哪些地方还未被探明、发生了哪些重要战役以及敌军在未来最可能在哪里发起攻击。影响力地图的结构还使得根据环境中不同位置的特征做出智能的推断成为可能。影响力地图能够挑选战略要地、查明敌方防御的薄弱环节、确定最佳的扎营点或易被攻击的区域、找出地形中的关隘以及确定玩家根据直觉或经验将选择的其他要素 (feature)。

并不存在惟一标准的创建影响力地图的算法，也不存在应用这种技术的惟一方式。本文将介绍多种比较流行的影响力地图概念，但这只是一个起点。创建和应用影响力地图的方式在很大程度上取决于游戏的战略和战术需求以及 AI 主体所处的游戏世界的设计。

### 3.6.2 一个简单的影响力地图

影响力地图 (Influence map) 几乎可用于任何类型的游戏世界地形中——方形网格、六角形网格或三维环境。出于简化的目的, 本文假设为 2D 网格, 这适用于大多数策略游戏; 而本文的最后一节将讨论影响力地图在更复杂环境中的应用。

我们首先将游戏世界划分为一系列的方形单元格, 所有单元格的初始值皆为 0。对于每个单元格, 加上某种类型的“影响力”。就这个例子而言, 我们将估算“战斗影响力”。因此, 对于每个友军个体, 加上一个正值; 而对于每个敌军个体, 则加上一个负值。

加上或减去的是个体战斗影响力的估算值。就现在而言, 我们假设每个个体的影响力为 1, 如图 3.6.1 所示。

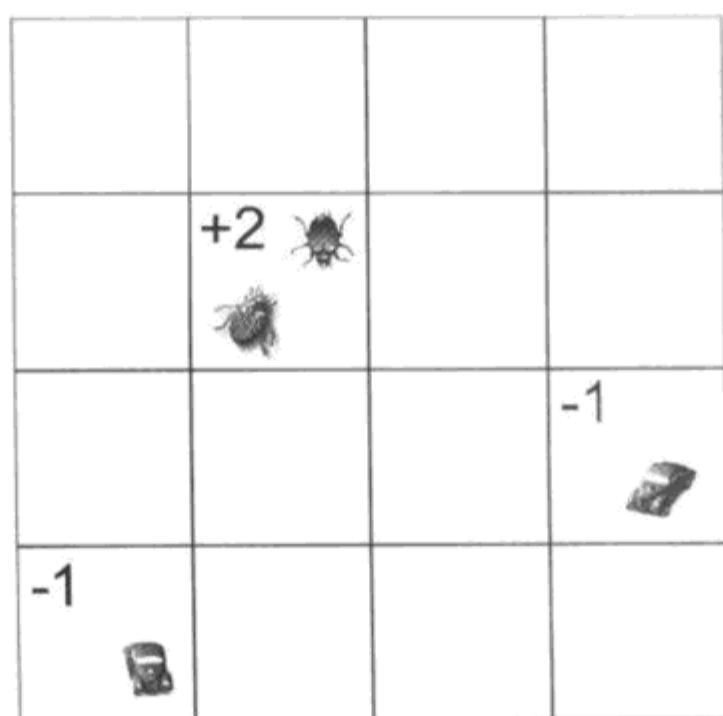


图 3.6.1 最初的影响力

接下来, 将每个单元格的影响力扩散到邻接单元格。假设单元格的影响力每向邻接单元格扩散一次, 其效果减半。因此, 如果单元格的值为 4, 则与之邻接的每个单元格的值将加上两点, 与之相距一个单元格的单元格将加上 1 点, 依此类推。

图 3.6.2A 说明了两个虫子的影响力是如何扩散的; 两个敌军个体 (福特汽车公司生产的穷凶极恶的主体) 的影响力以类似的方式传播 (图中没有画出), 但其影响力为负, 因为它们不讨人喜欢。

将车和虫子的影响力合并后, 将得到图 3.6.2B 所示的结果。这很好地指出了每个玩家发挥影响的区域。颜色较深的单元格属于我们, 而颜色较浅的单元格属于对方。更重要的是, 现在可以确定敌友之间的边界。如果两个相邻单元格的值分别为正和负, 则它们之间的网格线就是边界, 如图 3.6.2B 中的白线所示。

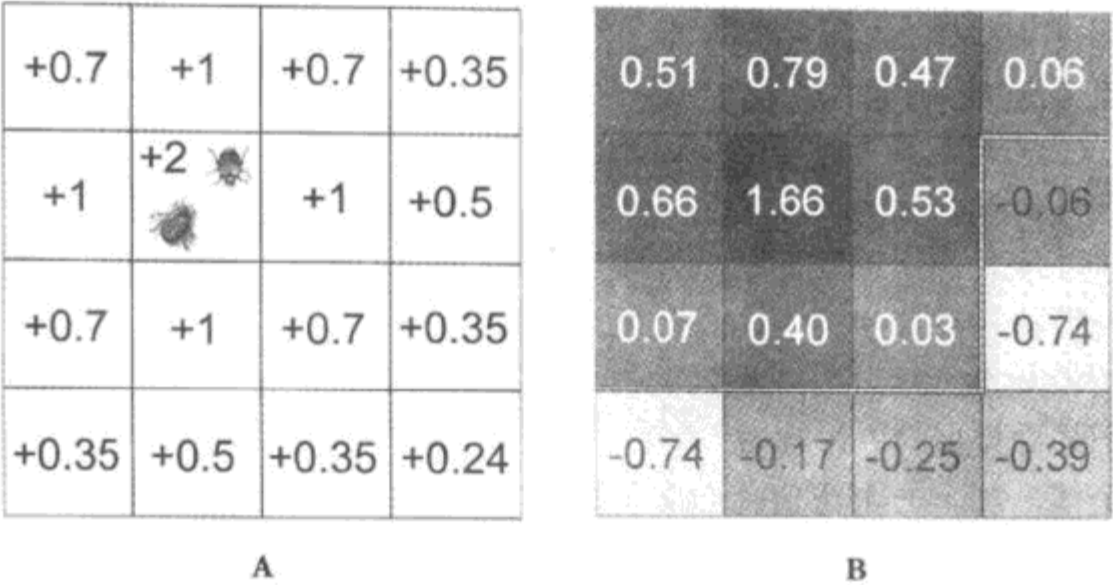


图 3.6.2 A)影响力的传播; B)最终的影响力地图

我们可以根据上述边界来决定如何部署兵力，以便进行进攻或防御。通过给敌军指定大于我军的权重（使用大于 1 的乘数），可以将边界收缩，摆出防御的姿态；如果给我军指定大于敌军的权重，则可以将边界向前推进，摆出进攻的姿态。

3.6.3 影响力地图中的单元格数据

显然，前面的例子用处不大。实际游戏中的影响力地图要复杂得多，每个单元格都存储了一些关于游戏世界的的数据，而不仅仅是一个数字。实际上，每个单元格都是一个小型数据库，其中存储了位于该单元格中的所有个体和资源的相关数据。下面是单元格中通常包含的一些统计数据类型。

- 战斗力 (combat strength): 单元格中当前个体的战斗力。应考虑诸如攻击力/防御力、生命值、攻击范围、火力 (rate of fire) 等因素。另外，可能还应根据游戏的设计，将个体划分为不同的类别，如远程打击还是近距离格斗，步兵还是骑兵，空军、陆军还是海军。
- 易被攻击的资产 (vulnerable asset): 单元格中玩家当前资产（如典型策略游戏中的村庄或军事基地）的估算值。
- 区域可见性 (area visibility): 一个数字，指出多长时间以来，该区域对玩家一直是可见或不可见的。
- 尸体数 (body count): 指出在该单元格中有多少个体死亡以及死亡的时间。
- 资源 (resource): 有待开采的各种资源，包括黄金、木材等。
- 通行性 (passablity): 通过单元格的难易程度，可按机动类型划分（飞越、步行、轨行）。这个值用于更准确地将单元格的影响力传播到邻接单元格，可作为特定决策的考虑因素。一种不错的方法是，分别存储 8 个通行性值，每个对应一个离开该单元格的方向。

影响力地图通常分别跟踪游戏中每个玩家的这些变量。可将此视为维护多个并行的影响力地图：每个玩家更新有关其资产的影响力地图以及表示其对其他每个玩家的了解情况的影响力地图。这很有用，因为它让您能够判断特定对手的强项和弱点，或在必要时将友军和敌军的影响力合并。然而，需要注意的是，当 AI 玩家超过 3~4 个后，性能将难以控制。

当然，也可以只为每个玩家维护一个影响力地图，并允许所有的 AI 玩家访问它。在有

隐藏地图 (hidden map) 或战争烟雾 (FOW) 的游戏中, 这构成了“作弊”, 在某些情况下可能导致次优的行为。

### 3.6.4 计算合意值

与其直接将单元格的基本统计数据作为决策根据, 不如将它们合并为一个“合意值 (desirability value)”。这是一个通过计算得到的值, 它指出了单元格对特定决策的合意程度。通过比较不同单元格的合意值, 可以确定对某种任务而言, 哪些单元格比其他单元格更合适。

最有用的计算合意值的公式常常是简单的加权和。根据要作的决策, 选择每个单元格中与之相关的变量, 并根据每个变量对决策的影响程度乘以一个相应的系数, 然后将乘积相加, 便可得到合意值。

计算不同的合意值时, 具体选择哪些参数在很大程度上取决于游戏的需求和游戏设计的独特特征。系数的选择也是主观的, 必须通过仔细的调整才能获得最佳结果。模拟退火 (simulated annealing) 或竞争进化 (competitive evolutionary) 方法是可行的, 但可能不可取。需要注意的是, 您可能需要对用于统计数据的不同度量单位 (如个体的生命值、火力、攻击力等) 进行补偿。

下面是一些合意值:

- 攻击和防御愿望 (attack and defense desirability)。我们通常可以分别计算玩家及其敌人的“脆弱”分数。敌人的脆弱分数高表明我们可以轻松地破坏该区域中的资产, 因此应考虑攻击该单元格中的敌人; AI 玩家的脆弱分数高表明我们在该单元格中有重要的资产, 易被攻击, 因此应更小心地保护它们。

通常, 有很多资产和重要资源但附近的守军很少的区域脆弱。因此, 如果对方玩家在某个单元格中的资产值为 80 (表示其基地建筑和资源), 而进攻值为 60 (表示守卫的敌军), 则脆弱值为 20。

- 探索 (exploration)。对于使用隐藏地图或 FOW 的策略游戏而言, 优秀的 AI 玩家将经常派遣其侦察员, 以刷新其战场视图。一种不错的探索方法是, 优先探索不可见时间最长的影响力地图单元格。进行这种决策时, 其他重要的因素有单元格中的敌军影响力和区域的通行性 (passability) (以便侦察员遭到攻击时能够逃跑)。

- 防御资产的位置 (defensive asset placement): 不能移动的防御资产应部署在脆弱资产的附近。它们必须位于脆弱性足够高因而有必要进行保卫的区域, 但这些区域的脆弱性又不能太高, 以至于无法制造它们。

关隘也是部署防御资产的好地方。通过预先计算通行性值, 可以很容易地确定影响力地图上的地形关隘点: 关隘点是影响力地图中这样的高通行性单元格, 即与其他高通行性区域相连, 但周围是通行性低的单元格。

- 资源采集资产的位置 (resource-collection asset placement): 位于离大量可开采资源最近的易于防守的区域时, 用作资源采集点的资产的效率将最高。

- 设备生产资产的位置 (unit-producing asset placement): 设备生产资产 (如兵工厂) 通常应部署在离敌军最近的防御区域中。

- 脆弱资产的位置 (vulnerable asset placement): 需要保护的资产应部署在防守最严密



的区域，并尽可能远离潜在的威胁。通常，将这些资产部署在不容易进入的区域，使其免受攻击也是一个不错的主意。对于扁平的长方形游戏世界，在地图的角上通常进入的通道较少，脆弱性较低，因此可以将地图边角的合意值设置得较大。

### 3.6.5 确定最佳的单元格大小

从某种程度上说，影响力地图的单元格大小可以是任意的。如果单元格太大，则影响力地图将难以识别诸如地形关隘点或敌军防御弱点等小型要素；如果单元格太小，则难以控制，导致大量的重复计算，并可能使用大量的内存。

在实践中，最好适当地设置单元格的大小。不要想当然地认为单元格越小，AI 越聪明。就典型的策略游戏而言，建议将单元格的宽和高设置为刚好容纳 10~20 个标准“个体”并肩排列，然后仔细地调整单元格大小，以获得最佳效果。

有些读者可能注意到了，在地图上随意地放置单元格可能导致问题。如果个体跨越两个影响力地图单元格，则将其影响力放置在不同的单元格中将导致不同的结果。由于下一节将描述的影响力传播，这通常不是什么问题。然而，一种处理这种问题的不错的方式是，定期地（可能是每次重新计算影响力地图时）调整整个影响力地图的空间偏移量( $X, Y$ )，使用随机或周期性的偏移值。这类似于在海洋中飘动的渔网，由于海浪的冲击而前后移动。

### 3.6.6 影响力传播

计算影响力地图中每个单元格的初始值后，接下来需要将每个单元格的值传播到附近的单元格，就像前面的例子中那样。这一过程也被称为修匀 (smoothing) 或模糊化 (blurring)，因为它与标准的 2D 图像模糊化技术有很多共同之处 (参见[Evans01])。

影响力传播提供了更精确的有关当前战术形势的图景。我们不但关心个体的位置及其当前行为，还关心它们可能做什么——它们可能影响哪些区域。如果我们的两门“射手”加农炮被大型 Plasma 坦克部队包围，我们希望 AI 知道，该区域实际上已被敌军占领。我们需要将坦克的影响力传播到“射手”加农炮所在的单元格。

传播只不过是使用“衰减规则”将每个单元格的影响力扩散到邻接单元格，衰减规则决定了给定单元格的影响力在地图上扩散时，将如何随距离的增大而降低。选择何种衰减规则有一定的主观性，并不存在惟一可行的技术——同样，需要进行调整，以获得最佳结果。作者发现，指数衰减最有用：选择一个 0~1 的衰减常数（通常为 0.6~0.8），然后每次将影响力扩散到邻接单元格时，将该常数作为乘数。假设衰减常数为 0.75 (75%)，则邻接单元格的值将为原值的 75%；与之相距一个单元格的单元格的值为原值的 56% ( $0.75^2$ )，依次类推。衰减常数必须与单元格大小相称：要将影响力传播相同的距离，则单元格越小，需要的衰减值将越大。

其他有用的衰减规则有线性衰减（每次扩散到邻接单元格时，影响力减小固定的值）和高斯过滤器（参见[Evans01]）。

注意，如果使用浮点数，则不管传播距离多远，传播的影响力将永远不会为 0。这意味着每个单元格都会将其影响力扩散到影响力地图中的所有单元格。所有的文献都认为这种现



象是件“糟糕的事情”。对此，最简单的解决方案是，当影响力值小到一定程度（通常当修匀后的影响力太小，不会带来任何差别）时，终止传播。这个终止常数最好通过实验来确定。

然而，请注意让影响力传播相当大的距离通常是个不错的主意。如果影响力地图包含大量的小单元格，而单元格影响力的传播距离不大，则影响力地图中很可能有大量的空单元格（即其值为 0），在这种情况下将难以准确地确定边界的位置。使用大型单元格，并将影响力传播较大的距离。

还有另一种影响力传播技术，它是基于四叉树的。所有单元格的值都沿四叉树传递到上层，这样可使用高层的四叉树单元格来获得其子单元格的大概“修匀”值。不幸的是，从某种程度上说，这种方法以随意的方式传播单元格的影响力。单元格影响力的传播距离与四叉树的结构密切相关，同时影响力沿某些方向的传播距离可能比其他方向远得多。作者发现，与前面介绍的传播技术相比，这种技术的灵活性和准确性都较低。

3.6.7 考虑地形

这里描述的传播技术并非在所有情形下都能绘制准确的图景。假设强大的敌军巫师在山脉的一边有一个炮台，这种传播技术将跨越山脉传播炮台的影响力，虽然巫师不能跨越山脉攻击我们，也无法调遣军队将我们包围。

要考虑地形对影响力地图（Influence map）的影响，方式有多种。也许最简单的方式是使用预先计算的每个单元格的通行性值，将其作为衰减值乘数，如图 3.6.3 所示。每个单元格要么包含一个通行性值，要么包含 4 个或 8 个通行性值（对应于离开单元格的主要方向）。然后以类似于宽度优先搜索算法或淹没填充（flood-fill）算法的方式将影响力从单元格传播出去。虽然图 3.6.3 中没有指出，这也能够处理影响力已减小但未完全阻断的单元格。

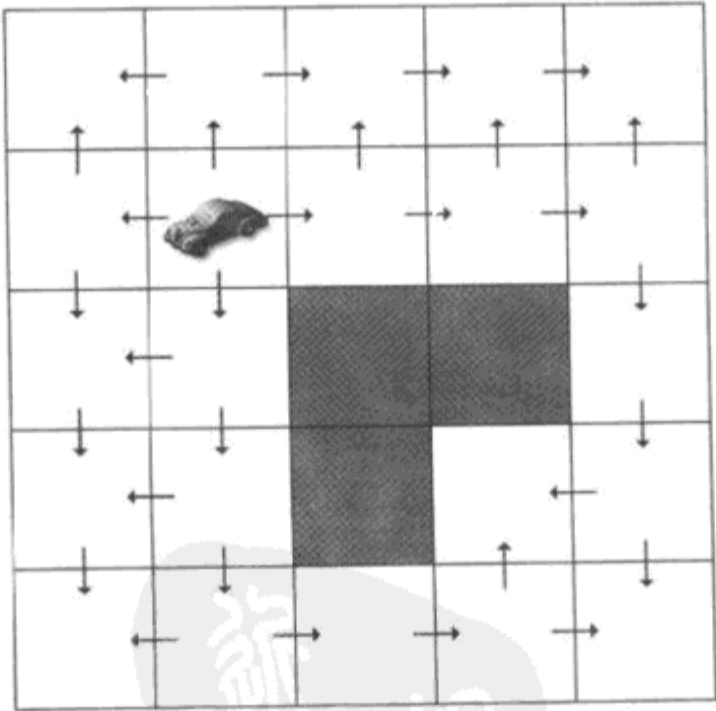


图 3.6.3 根据地形进行传播

另一种技术需要预先计算邻接单元格之间所有可能的路径（如图 3.6.4 所示）。对于每个单元格，我们在地图预处理期间执行寻径步骤，计算从该单元格到附近所有单元格的最短路

径(如果最短路径小于最大路径距离)。然后将计算得到的到每个目标单元格的距离存储起来,并在传播影响力时将其作为到附近单元格的实际距离。如果没有到某单元格的路径,则认为影响力无法传播到该单元格。

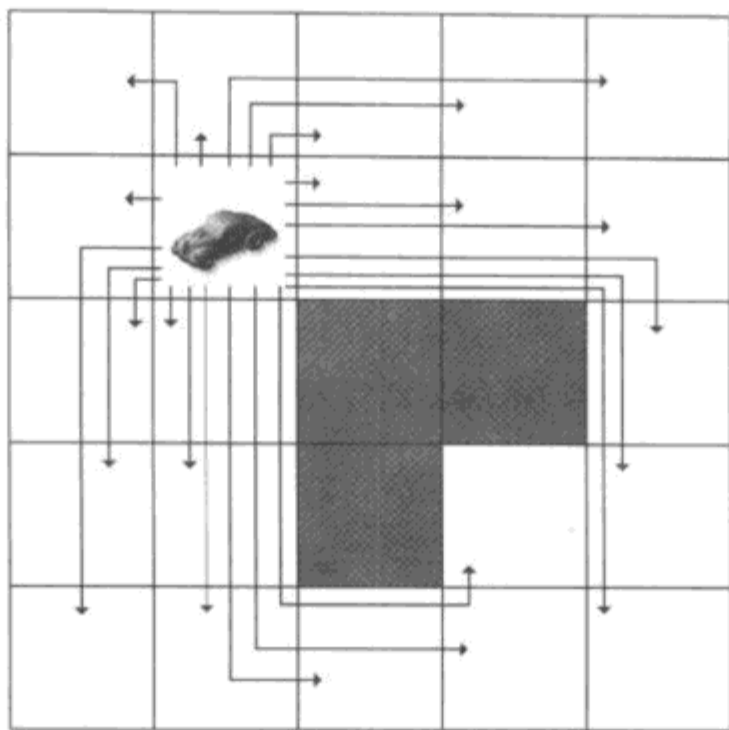


图 3.6.4 预先计算的传播

这种技术的优点是,提供了准确的影响力传播。如果存在跨越山脉的简易途径,如穿越山脉中心,则这种影响力地图传播将准确地反映这样一个事实,即山脉并非重要的战术屏障。

不幸的是,这种技术难以应用于动态环境中。如果游戏世界允许玩家建造长长的城墙或封锁山脉通道,则预先计算的传播值反映的将不再是真实情况,而实时地更新影响力地图是极其困难的。这种方法还可能需要大量的预处理时间,因为它要求找出影响力地图中每个单元格到其他单元格(可能有几十个,甚至几百个)的路径。

### 3.6.8 特别考虑

地形的影响随军种而异。山脉不能阻止空军前行,而海军的影响力只能在水上传播,不能在陆上传播。因此,每个影响力地图单元格必须分别计算不同的军种(例如,每个单元格中的“飞行资产”和“非飞行资产”),并根据地形以不同的方式传播它们的值,这很重要。

有些军中的火力(rate of fire)范围很远,同时,如果个体的火力范围为特定的值,则在它可以移动到的任何位置,其火力范围也应与此相同。要针对每个影响力地图单元格考虑这一点,一种不错的方式是,按其火力范围分别跟踪远程打击个体(ranged distance unit)(可能按影响力单元格宽度的倍数进行分类)。通过传播扩散这些远程打击类别的影响力,然后将其影响力从受影响的每个单元格再扩散  $N$  个单元格,且影响力值保持不变。这样战舰便能够打击远离海滨的内陆地区,虽然它无法登陆。

您可能还会发现,根据未来的位置而不是当前位置将移动装置添加到地图中很有用。这样,影响力地图将更准确,尤其当不频繁地重新计算影响力地图时。最简单的方法是航位推测法

——根据其速度向量预测 5~10 秒后的位置。由于您可能要编写 AI，因此如果个体要进行寻径，您也需要查看每个个体的未来位置（虽然如果 AI 玩家查看其他玩家选择的路径，则有作弊之嫌）。

### 3.6.9 刷新影响力地图

如果您的 AI 需要定期分析影响力地图的很大一部分，则应定期地（也许是每隔 1~10 秒）重新计算整个影响力地图。鉴于典型策略游戏的步伐，快速的刷新可能并不能使 AI 更有效。计算好影响力地图后，便可以使用它来快速地完成很多计算。

另一种方法是按需刷新，这是一种延迟计算（Lazy evaluation）技术。这种方法更灵活，同时如果不需要进行大量的影响力地图分析，则其效率也更高。采用这种方法时，仅当给定的单元格被查询时才计算其中的值——搜索位于最大距离范围之内内的所有单元格，以了解它们的值如何传播给原来的单元格。这种技术还有另一个优点——可以在查询时指定传播参数和合意值。

### 3.6.10 3D 环境中的影响力地图

前面重点介绍的是影响力地图在 2D 环境中的应用。然而，创建影响力地图及相关的方法也适用于更复杂的环境，如典型动作游戏中的 3D 环境。

将 2D 网格用于创建 3D 影响力地图是个糟糕的主意，因为这将不能准确的反映游戏环境的地貌。幸运的是，3D 环境中的 AI 寻径通常（虽然并非总是）是基于一种导航网（navigation mesh）方法的（参见[Snook00]）。导航网由一系列互联的凸多边形组成，它描述了角色在游戏世界的哪些地方可以移动。我们可以将导航网中的每个多边形作为一个影响力地图单元格。多边形之间的连接描述了影响力的传播途径。影响力将随其沿每个网节点传输的距离而降低。

由于 3D 游戏世界的地貌比 2D 世界要复杂，关注的重点是各个战士，因此地形评估通常比影响力地图提供的玩家对玩家的实时战术评估更重要。AI 主体必须能够辨别不同区域的战术意义，这至关重要。下面的列表解释了一些战术评估因素：

- **易被攻击性（vulnerability）**：3D 动作游戏通常涉及火力（rate of fire）强大的远距离打击武器，因此考虑影响力地图中每个单元格的火力覆盖范围至关重要。AI 主体常常需要判读给定单元格的“覆盖”程度。一种简单的方法是，计算从每个影响力地图节点发射时，立方体的每个面的覆盖程度。

然而，我们经常需要知道给定的单元格（可能的目的地）是否在另一个单元格（敌人的位置）的射击范围内，或者从前者是否能够射击到后者中的目标。优秀的发射线表示法将列出从给定节点能够攻击到的所有节点，但对于大型的高度互联的环境（其互联度为  $N^2$ ），这种方法将失控。一种解决方案是，将影响力地图节点划分为“区（zone）”，这样位于给定区中的所有节点都在同一个“房间（room）”中，然后使用这种表示法来判读哪些区在其他区的攻击范围内。

- **可见性（visibility）**。这与易被攻击性类型，只是它考虑不是武器的火力范围，而是光照水平，并能够穿过武器火力无法穿过的面，如防弹玻璃。使用动态灯光，且可被开

关时，这种计算将非常棘手。

- **通行性 (passability)**。与前面一样，这指的也是通过某个区域的难易程度。狭小的通道通常更难以通过，而电梯、梯子等通道导致移动缓慢，因此通行性较低。
- **高度优势 (height advantage)**。无论是进攻还是防御，海拔高度比周围高的位置都有战术优势，尤其是当游戏中可以使用手榴弹时。

这些是影响力地图单元格中需要预先计算的基本统计数据。现在，我们使用一些系数，计算合意值。

对进攻而言，通行性高、覆盖范围广、可见性低，同时对很多高可见性（最好且通行性低）的区域有良好的发生线的位置是最佳的。好的进攻位置附近通常还应有良好的掩体位置，以免主体受到攻击。

对于防御而言，这样的位置是最佳的，即覆盖范围好、可见性低、且可能对其发起攻击的所有位置的可见性都很高。

另外，如果预先为所有节点计算这些合意值，然后便可以使用标准的影响力地图创建技术将这些值传播给邻居，最后获得完整的战术评估。

最后请注意，也可以使用影响力地图来推断对手的情况——可以判断其覆盖程度、可见性、通行性和高度优势。在任何时候，可以根据这些信息选择最适合向其发起攻击的对手；或搜索一个目的地，在那里我们比所有的对手都有优势。

另一个有用的扩展是，估计对手在不远的未来最可能位于什么位置，这是通过找出最适合他们的目的地来实现的。这让我们能够针对敌人的行动做好准备，预先埋下伏兵。

### 3.6.11 参考文献和推荐读物

[comp.ai.games95] The seminal 1995 “influence mapping” thread on comp.ai.games (various authors) is reprinted at [www.gameai.com/influ.thread.html](http://www.gameai.com/influ.thread.html).

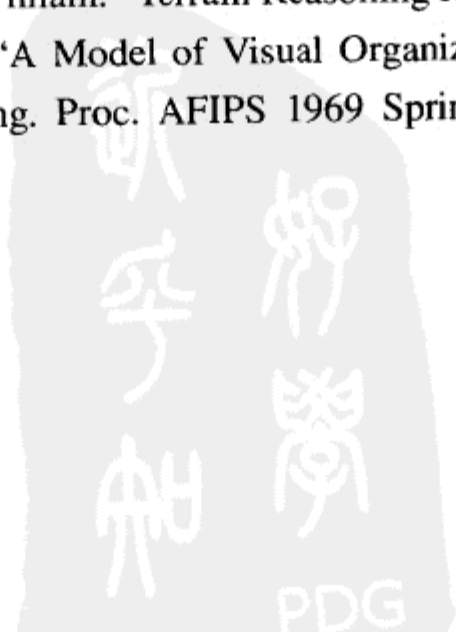
[Evans01] Evans, Alex. “Four Tricks for Fast Blurring in Software and Hardware,” available online at [www.gamasutra.com/features/20010209/evans\\_01.htm](http://www.gamasutra.com/features/20010209/evans_01.htm).

[Pottinger00] Pottinger, Dave. “Terrain Analysis in Realtime Strategy Games,” available online at [www.gdconf.com/archives/proceedings/2000/pottinger.doc](http://www.gdconf.com/archives/proceedings/2000/pottinger.doc).

[Snook00] Snook, Greg. “Simplified 3D Movement and Pathfinding Using Navigation Meshes” from *Game Programming Gems I* (Ed. Mark DeLoura, Charles River Media, 2000).

[Sterren01] van der Sterren, William. “Terrain Reasoning for 3D Action Games.”

[Zobrist69] Zobrist, Albert L. “A Model of Visual Organization for the Game of Go.” The seminal article on influence mapping. Proc. AFIPS 1969 Spring Joint Computer Conf. 34: pp. 103~112.



## 3.7 策略评估技术

Paul Tozour

gehn29@yahoo.com

本文介绍一套对策略决策很有帮助的技术。前一篇文章介绍了如何创建影响力地图，提供了一种地理级战术评估方式，这些数据结构为AI主体或玩家提供了在战略和功能级对游戏的当前状态进行评估的方式。

这些技术适合于涉及到经济管理、资源分配决策和/或技术进步的游戏中，如策略游戏、经济状况和“神话游戏（god game）”。然而，它们也可用于其他游戏类型中。

### 3.7.1 资源分配树

资源分配树（resource allocation tree）是一种树结构，表示玩家控制的所有资产的特定用途。这种树以层次结构的方式，按功能对当前所有的个体和资源进行分类。

这种表示法的用途在于，它让AI玩家能够对游戏中所有玩家的战略优点和弱点进行评估，还为各种经济生产和资源分配决策提供了很好的基础。例如，它为生产什么样的设备以及如何重新将已有的设备指定为不同的功能角色提供了直接的决策依据。

树的最顶端是根节点，表示玩家的所有资产。根节点的下面是游戏中可用资产的主要类别。

假设主要类别有军事、经济 and 智能。其中每种类别被进一步划分为各种子类别。例如，军事可能分为进攻和防御，而这些子类别还可以划分为弹道式（Ballistic）、远程（Ranged）、格斗（Melee），以指出军事设备的功能角色。经济主类别可划分为资源采集、设备生产、基地建设、技术进步等。图3.7.1是资源分配树的一部分。

树叶为个体的具体类型。例如，长矛兵的位置为根/军事/进攻/格斗/长矛兵。“长矛兵”节点本身可能包含关于所有长矛兵的统计数据（如总数、攻击力、战斗中牺牲的长矛兵数等），以指出其进攻能力。

和影响力地图一样，AI玩家也应针对游戏中的每个玩家，维护一个这种数据结构的实例。其他玩家的图形表示当前玩家对每个玩家的策略资产的了解情况和功能划分的估算。

长矛兵主要是防御个体，但也可用于攻击或探索。这提出了一个问题，对于能够扮演多种功能角色的个体，如何对其进行分类。



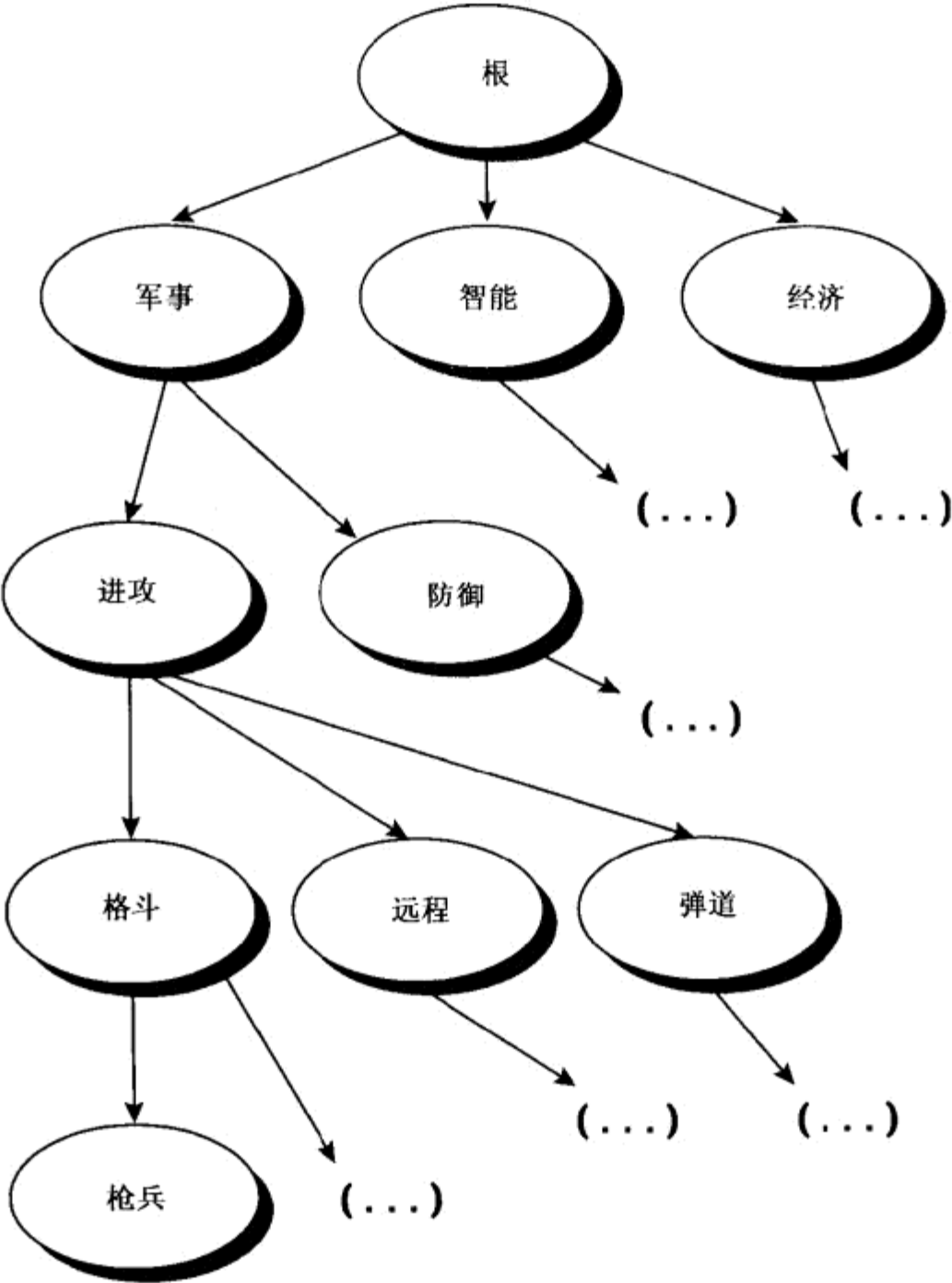


图 3.7.1 资源分配树

在任何给定的时刻，AI 玩家通常只给给定个体指定一种功能角色，因此作者建议根据个体当前的职业对其进行分类，而不要将单个长矛兵分为 10%的探索、60%的防御和 30%的进攻。位于进攻、防御和探索下面的“长矛兵”节点应表示当前被指定为相应功能角色的长矛兵。

另外请注意，在任何时刻，资源分配树中都只包含我实际拥有或能够生产的资产，这很重要。如果我没有魔法师堡垒（Mage Tower），且没有能力建造这样的堡垒，则没有理由将其作为树中的一个节点。

### 3.7.2 计算希望的资源分配

资源分配树（Resource allocation tree）结构提供了一种简单、自然的确定合适资源分配的方式。可以从树根开始依次向下进行处理，给树中的每个节点指定优先级。

首先，将根节点的值指定为 1.0，这表明希望的资源分配比例为 100%。在每个节点，将



其当前值按一定比例分配给各个子节点。从根节点的值 1.0 开始，将它分为军事 0.3、智能 0.1 和经济 0.6。在“军事”节点，将其值的 67%和 33%分别分给防御和进攻（防御为 0.2，进攻为 0.1）。沿树向下不断重复这一过程，为每个树节点指定“希望的分配”值。

划分每个节点值的算法取决于游戏的设计。为响应游戏世界的状态变化，每个非叶节点都必须有定制的逻辑代码，以不断更新其值在子节点之间的分配。这需要不断进行调整，以获得最佳结果。通常，这样做是个不错的主意，即您对哪些因素将导致计算机玩家修改权重有深入了解之前，进行一些猜测并使用预先确定的常量。

### 3.7.3 判断当前的分配情况

同时，我们可以使用资源分配树来计算每个节点的当前分配值。这让我们能够将实际拥有的资产进行划分。

这种计算的方向相反——从最下面开始，沿树向上到达根节点。例如，我们遍历当前负责防御的每个长矛兵，并计算一些值——它指出每个防御长矛兵的估算值（如将长矛兵的攻击力同其生命值相乘）。然后将这些长矛兵的得分相加，并将结果作为防御/格斗下面的“长矛兵”节点的当前分配值。接下来，将这个值向上传递给父节点（格斗），后者将其从子节点那里收到的所有数字相加，并将结果传递给它的父节点（这里为“防御”节点）。最后，根节点接收其所有子节点的输入，并将它们相加，得到一个很大的浮点数，这个数指出了我们当前拥有的所有资产。

计算出“总资产”值后，便可以将各个节点的值重新规格化到 0~1 之间。方法是，再次访问树中的每个节点，将其“当前分配”值除以根节点的值，这样，根节点的值将为 1.0。

至此，我们可以直接将每个树节点的当前分配值同期望的分配值进行比较，以判断超过（或少于）预算多少。

### 3.7.4 策略决策

显然，最终的分配树表示提供了一种优秀的、维护军力平衡的方式。如果我在一场战役中失去了所有的 20 头大象，则“大象”节点及其所有父节点的分配值将小于希望的分配值。

当然，这并不意味着我一定要用新的大象替代所有的大象。当“大象”节点的父节点获得填充新的子节点所需的资源后，它可能认为最佳的措施是制造一个 Plasma 坦克营，因为现在已获悉了制造这种坦克所需的技术。

资源分配树（resource allocation tree）主要用于判断制造什么样的新设备以及如何给已有的设备指定合适的角色。通常，首先需要找出哪些迫切需要增加资源的节点，然后确定为解决这种问题，是对已有个体重新布防合适还是生产新个体合适。

资源分配树还为设计独特的玩家个性提供了很好的途径。无论开发的是扩张主义者、穷兵黩武的“成吉思汗”、关心经济的资本家还是强调研究的技术政论家，都只不过是调整资源分配树中相应部分的系数，以侧重或冷落特定的节点而已。通过稍微调整资源分配树的不同部分，可使 AI 玩家侧重于某种个体，在经济增长和国防之间取得不同的平衡，喜好特定的战略资产类别，甚至改变 AI 玩家的风格。

当关心的是战斗时,这样做常常是个不错的主意,即维护一个预先计算的“战斗平衡表”(其中包含个体的相对强弱),并将其用于资源分配树的“军事”分支中的决策。这是一个 2D 查找表,让您能够判断任何个体同其他个体战斗的效果。通过分析表示您对特定敌人的了解情况的功能资产树,可以知道敌人的军力组成情况,并重点生产能够最有效地打击敌人的资产。

最后,资源分配树适合存储其他的各类统计数据。可在影响力地图单元格中存储的几种因素也可在功能资产树节点中存储。

另外,记录树中的哪些节点被证明是有效的以及哪些节点遭到敌人的攻击也是一个不错的主意。前者让我知道长矛兵在对付另一个玩家时很有效,因此我将资源更多地分配给让长矛兵取得胜利的功能角色;后者让我知道敌人倾向于攻击我的资源采集设备,因此我在以后将采取更多的保护措施。

### 3.7.5 值的估量

采用这种技术时,面临的最重要的挑战可能是找出合适的方式来估量每个个体在其特定树分支中的值。每个节点中使用的数值必须与其分支中所有其他个体的值相称。“军事”分支中的值表示个体对战斗的贡献,需要考虑攻击力、发射速度、移动速度、防护值(armor value)、当前的生命值以及其他参数。“智能”分支中的值需要考虑与个体的探索能力相关的因素,如望距离和移动速度,但无需考虑攻击力或生命值。对于“资源采集”节点(位于“经济”分支中),应考虑资源采集设备采集资源、卸载资源以及再次返回资源区的速度。和其他游戏 AI 一样,要找到关联这些值的合适方式,需要进行试验,以获得对您的游戏而言最佳的解决方案。

另一个可能的挑战是处理有多种资源的经济体系。在生产不同个体需要耗费不同数量的黄金、能源和 Tiberium 的游戏中,单个分配值并不对应一种资源。这里不打算介绍这种问题的解决方案。

### 3.7.6 依存图

依存图(dependency graph)是一种数据结构,指出了游戏中不同资产类型之间的所有依存性。依存图包含所有基于依存的关系,如游戏的“技术树”和“建造树”。

主要的依存类型是创建(creational)依存。它指出了要建造特定的资产必须满足的一些条件。例如,要创建长矛兵,必须有兵工厂(Barrack)。要进入帝国时代(imperial age),必须建造一个城堡(Castle)。

创建依存还可以包含资源依存和其他更抽象的依存。兵工厂需要黄金和木材,而黄金和木材来自农民的劳动。农民从金矿中采黄金,从森林中砍木材。

图 3.7.2 是一个小型的依存图,其中只包含创建依存。农民能够建造兵工厂和射箭场,但只有进入中世纪时代后,才能建造射箭场。

第二种依存是支持依存。魔法师需要神力,而后者只能来自神殿。没有神殿,魔法师将毫无用处,因为没有神殿提供宝贵的神力,他将无法制造符咒。

和前面讨论的数据结构一样，AI 玩家应维护多个并行的依存图，自己和其他每个玩家各一个。

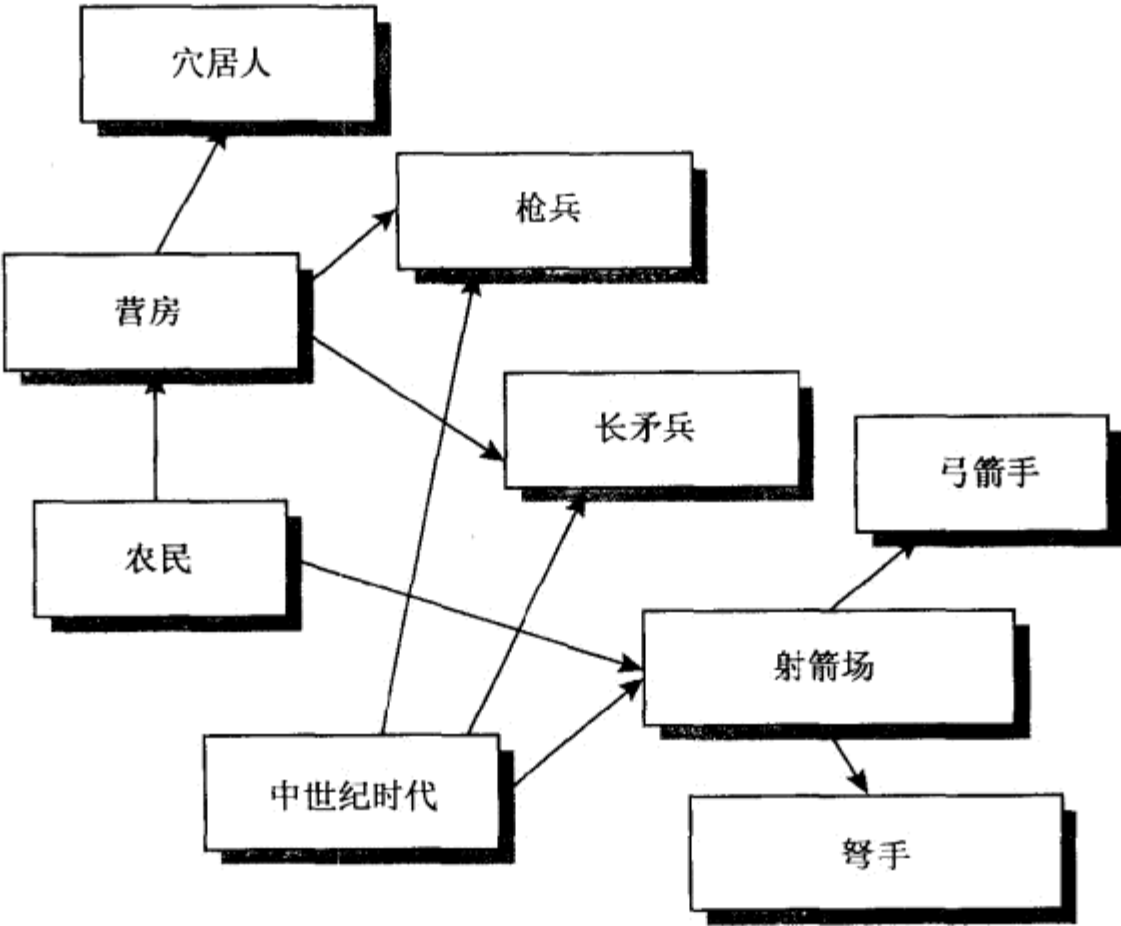


图 3.7.2 依存图

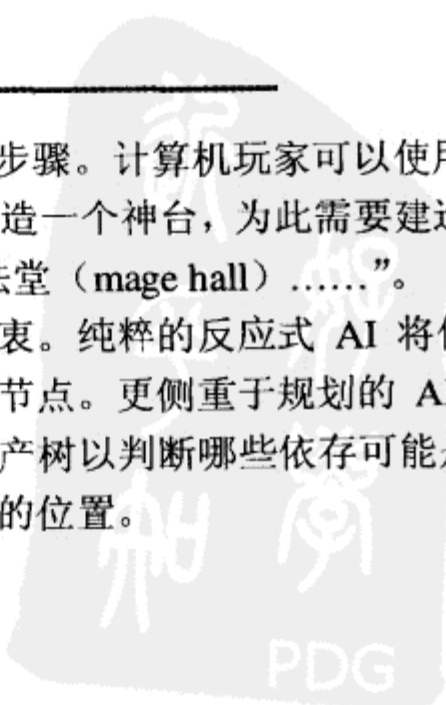
3.7.7 依存图节点

依存图 (dependency graph) 中的一个节点通常包含多种不同类型的数据。其中有用的类别有：玩家当前拥有的（如果考虑的是另一个玩家，则认为是它拥有的）某种类型的个体总数、这些个体的总估算值以及当前正在生产的这种个体（例如，由兵工厂生产）的总数。由于关心的内容有交叉之处，因此依存图中的节点可能类似于功能资产树中的节点，甚至它们的物理数据结构相同。它们之间的主要区别在于，资源分配树只记录当前可用的资产，而依存图记录所有可能的资产。

3.7.8 经济规划

依存图最主要、最显而易见的用途是确定完成目标的步骤。计算机玩家可以使用依存图来确定为生产特定的资产，需要建造哪些东西。“我需要建造一个神台，为此需要建造一个藏经阁 (Library)，然后获得合适的神力源，再开始建造魔法堂 (mage hall) ……”。

首先满足哪种依存是着眼于当前和规划未来之间的折衷。纯粹的反应式 AI 将使用资源分配树列出它能够立刻生产的所有资产，并选择最合适的节点。更侧重于规划的 AI 将分析图中的所有节点以确定值得追求的长远目标，查询功能资产树以判断哪些依存可能是最有价值的，并向最有前途的技术迈进，而不管它位于图中多深的位置。



诸如，当满足依存的方式有多种时（通过 A 和 B，都能到达 C），这种处理方式将变得微妙起来。在实践中，这种情况很少发生，因为游戏设计者聪明地避开了这种类型的依存。当可选的依存很多，以至于到给定节点的最佳路径并非显而易见时，可以使用任何标准的搜索算法来快速地解决问题。

### 3.7.9 查找脆弱的依存

依存图可用于分析玩家的优势和弱势以及查明对方最脆弱的依存。在《红色警报 II》中，敌方 AI 通常会在摧毁我的战争工厂（War Factory）后破坏我的兵工厂。这是一种聪明的破坏性策略，因为这迫使我花费宝贵的时间和金钱重建兵工厂，然后才能重建战争工厂。

判断依存图（dependency graph）中给定节点是否脆弱的因素有三个：

- **内在价值（intrinsic value）**：有些资源本身很宝贵。核武器发射井很宝贵，因为能直接攻击敌人。资产在图中的位置越深（沿技术树的路径越长），通常其内在价值越高。
- **强的子依存（strong child dependency）**：有些资产会由于其能够创建或支持的东西而值得攻击。战争工厂能够制造坦克和其他车辆。藏经阁让我们能够建造魔法堂，进而得到魔法师。核电站（fusion plant）为玩家的基地提供电力（一种支持依存）。
- **弱的父依存（weak parent dependency）**：我们也可以破坏父依存，就像在摧毁战争工厂后破坏兵工厂那样。父节点相对较弱（支持资产不多）且易于攻击（保护不好）的节点容易受到这种攻击。

无论是攻击还是防御，我们都可以采用相同的推理。对于攻击，我们通常选择对方最宝贵的依存，首先攻击战争工厂，然后是兵工厂。防御时，我们根据依存图来提防这样的攻击——加强防御，并建造备用的建筑。

### 3.7.10 策略推理

依存图的一个优点是，为根据不完全的资料推断其他玩家的当前资产和可能采取的策略提供了依据。例如，如果我知道敌人有一个兵工厂，则可以确信他要么已经有长矛兵，要么不久就能生产长矛兵。同样，如果我发现了敌人的长矛兵，则可以 100%地肯定他在附近有兵工厂（至少，在生产长矛兵时有——生产长矛兵后，他很可能摧毁了兵工厂）。

推理方向有两个：向前推理和向后推理。

向前推理时，我们知道特定的玩家拥有特定的个体或资源，因此断定他很可能满足子依存。见到兵工厂后，我们断定敌人有长矛兵。

向后推理时，我们沿依存链向后走，根据给定的个体断言（assert）使其依存将得到满足。见到长矛兵让我们确信存在兵工厂，虽然我们从来没有看到。

这种基于依存的推理可以进行很远。如果我看到敌军的大魔法师（Grand Mage），则可以断定敌人很可能有魔法师神秘建筑、高级魔法（High Magic）升级功能、神殿、圣人协会（Sage's Guild）、藏经阁以及让大魔法师得以出现的其他依存。然后，我可以对这些节点应用向前推理，推断出其他可能性。由于根据大魔法师可以推断出存在圣人协会，因此我认为该玩家很可能能够生产圣人。



这是一种概率推理，同被称为“Bayes network”的流行的推理技术极其类似（有关细节，请参阅参考文献）。

有趣的是，也可以根据推理断定某些节点不太可能出现。如果在游戏的某一点设置玩家经济状况的最大可能值——可用建造的时间或根据影响力地图中的数据得出的推理，则某些依存将导致其他依存得以满足的可能性较小。我知道，最优秀的玩家在4分钟内能够建造一个 Red Dragon Roost 或一个核武器发射井，但不能将两者都建造出来。因此，如果看到其中的一个，则另一个存在的可能性将较小。

当然，这里的工作量很大，但您可以作弊，直接查看其他玩家的资产。如何选择取决于您的良心以及您认为这种决策对游戏的娱乐性所产生的影响。

### 3.7.11 玩家个性

---

和功能资产树一样，我们可以使用依存图赋予 AI 玩家独特的个性。

优先需要调整的是依存图中各种节点的脆弱值。通过增大或缩小不同节点的脆弱性值，可以让 AI 玩家提高或降低这些资产的重要性。调整对手的依存图中的这些值将改变我们将这些敌方资产作为攻击目标的可能性；调整自己的依存图中的这些值将改变我们的经济发展方式和我们将采用的技术。

要设置初始优先级，一种简单的方法是选择给定 AI 玩家的一组终极“目标”。对于图中最右边（最深）的所有节点，找出一种合适的算法，根据合意性对这些依存进行排序，并最终得到良好的合意值。然后，将这些合意值传递到依存图的左边，从而明确地告诉 AI 玩家哪种技术更合适。

### 3.7.12 总结

---

AI 的很大一部分是决策，以良好的方式表示游戏状态将使 AI 玩家更容易进行决策。

本文和前一篇文章描述了这种决策的数据结构基础。虽然并非每个游戏都能使用所有这些数据结构，但无论您使用这些数据结构中的哪一种来进行战略和战术决策（tactical decision），在这种数据结构之间相互通信的机会都是无穷的。当影响力地图、资源分配树和依存图相互通信和共享其数据时，结果将是总体大于各个部分的和。影响力地图指出了敌人的位置，资源分配树指出了为攻击敌人需要哪些东西，而依存图指出了如何建造这些东西以及制造好后如何使其处于最佳状态。

### 3.7.13 参考文献

---

[Ferber99] Ferber, Jacques, *Multi-Agent Systems*, Addison-Wesley, London, England, 1999. A useful perspective is to consider the cells/nodes of the data structures described here as hierarchical “agents” in a multi-agent system.

[Pearl88] Pearl, Judea, *Probabilistic Reasoning in Intelligent Systems: Networks of Plausible Influence*, Morgan Kaufman, San Francisco, CA, 1988. A must-read introduction to Bayes nets and

probabilistic reasoning systems as they relate to AI.

[Russell95] Russell, Stuart and Norvig, Peter, *Artificial Intelligence: A Modern Approach*, Prentice Hall, Upper Saddle River, NJ, 1995. A comprehensive introduction to AI techniques—the “AI bible.”





## 3.8 3D 游戏中的地形推理

William van der Sterren, CGF-AI

william@cgf-ai.com

“位置，位置，还是位置！”这不但适用于房地产，也适用于 3D 动作游戏中的虚拟世界。在 3D 动作游戏中，位置与狙击点、据点、攻击方式、隘口和带“红色装甲”的供电区域一样至关重要。当位置在游戏中至关重要时，在游戏 AI 中也至关重要。

本文介绍对位置及其在游戏中的作用进行推理的技术，演示如何将位置概念转换为算法，以便 AI 能够计算并掌握这些概念。这有助于 AI 选择好的动作位置和了解其他行动者（actor）占据的位置，让 AI 处于有利的位置以帮助或挑战玩家。

首先，我们选择一种 AI 能够高效地处理的地形表示法：中继点图（waypoint graph）。为说明基于中继点的推理，将介绍一个范例问题。然后，确定战术特征并将其与中继点属性关联起来。我们将创建一些公式，使用诸如世界几何形状等静态数据以及诸如行动者行为等动态数据来计算这些中继点属性。最后，我们讨论如何将地形推理集成到游戏中，并介绍基于中继点的推理的其他应用。

### 3.8.1 以方便推理的方式表示地形

如果位置不多，则对其进行推断将很容易。然而，当前的游戏世界包含数以万计的、可进入的多边形。在游戏中，多个 AI 行动者将考虑可见或不可见的位置。如果以原始多边形的方式来完成这项工作，则消耗的资源将超过 3D 图形渲染。

更为复杂的是，位置的重要性并非完全由其特征决定，还取决于它与周边位置之间的关系。例如，是否容易进入？是否可从很多其他的位置看到？附近是否有 power up（威力升级）？另外，实际的游戏进展使问题更复杂：有些位置（如目标附近）常被玩家光顾，而其他区域不会。

因此，为有效地推断地形，必须使用一种近似的地形表示法，它不像原始多边形那样过多地纠缠于细节。这种表示法应易于表达位置之间的关系，支持捕获基于位置的游戏活动。最理想的情况是，这种表示法让我们能够使用预先计算的中间结果，以节省一些 CPU 时间，供更高级的 AI 使用或提高游戏的速度。

### 3.8.2 中继点 (waypoint)

---

大多数 3D 游戏 AI 采用一种易于处理的地形表示法。它们使用中继点、类似的导航辅助工具（如网格）或单元格（参见[Snooke00]、[Rabin00]、[Reece00]）。这些中继点表示地形中最重要的部分：游戏行动者能够进入的地形。每个中继点都是其周边环境的样本和近似。中继点数通常为 250~2500。

根据中继点进行推理很有吸引力，因为很多游戏 AI 已经使用这些中继点来移动、寻找路径、标记特殊的物品或障碍物以及接收来自关卡设计器（level designer）的提示。由于 AI 的很多行动都考虑了中继点，并认为玩家位于附近的中继点，因此捕获每个中继点的游戏数据很容易，几乎没有开销。

要对基于中继点的地形表示进行推理，这种表示必须粗略地说明地形及相关的属性。中继点网络的密度应足够高，能够指出所有的相关位置以及所有的掩体和隐蔽点。通常，这意味着需要的中继点数将超过 AI 导航和寻径的要求。

除最短路径和移动外，地形推理常常还需要处理中继点之间的其他关系。为说明为何需要其他关系以及如何推断中继点，最佳的方式是使用一个范例。

### 3.8.3 范例地形和 AI 需求

---

为说明基于中继点的推理，来看一个例子：我们希望 AI 行动者在战斗开始之前和战斗期间，从图 3.8.1 所示的区域中选取最佳的进攻和防守位置。

为支持 AI 有效地挑选并进入最佳的位置，对于每个中继点和各种不同的方向，我们计算相应的进攻值和防守值。这种值是根据中继点图和游戏世界几何形状计算得到的，并基于玩游戏的经验对其进行改进。这种对每个中继点的战术认识可作为寻径（pathfinding）、物群模拟（flocking）的输入，以选择一个将目标尽收眼底的守卫位置。

该范例区域有一个目标和两个入口，其中分布了一系列稠密的、以网格方式排列的中继点。中继点本身提供了良好的、关于关卡结构的线索（如图 3.8.1 中央所示）。

### 3.8.4 战术分析 (tactical analysis)

---

估算位置的战术价值时，需要考虑的因素很多。可将这些因素转换为中继点属性，然后对后者进行计算。我们从 *Quake* 式的夺旗（capture-the-flag）游戏和战术模拟的角度来考虑该范例区域。

在一个以快速移动、power up、非致命武器和火箭发射器为特征的快速夺旗游戏中，位置的战术价值在很大程度上取决于下述特征：

- 对攻击者而言，能够沿任何方向快速移动的位置很重要；
- power up 附近的位置颇有价值；
- 易遭受火箭袭击的位置不太好；
- 俯视通往旗帜的路径，且在向旗帜冲锋时不易发现的位置是非常好地防守点。

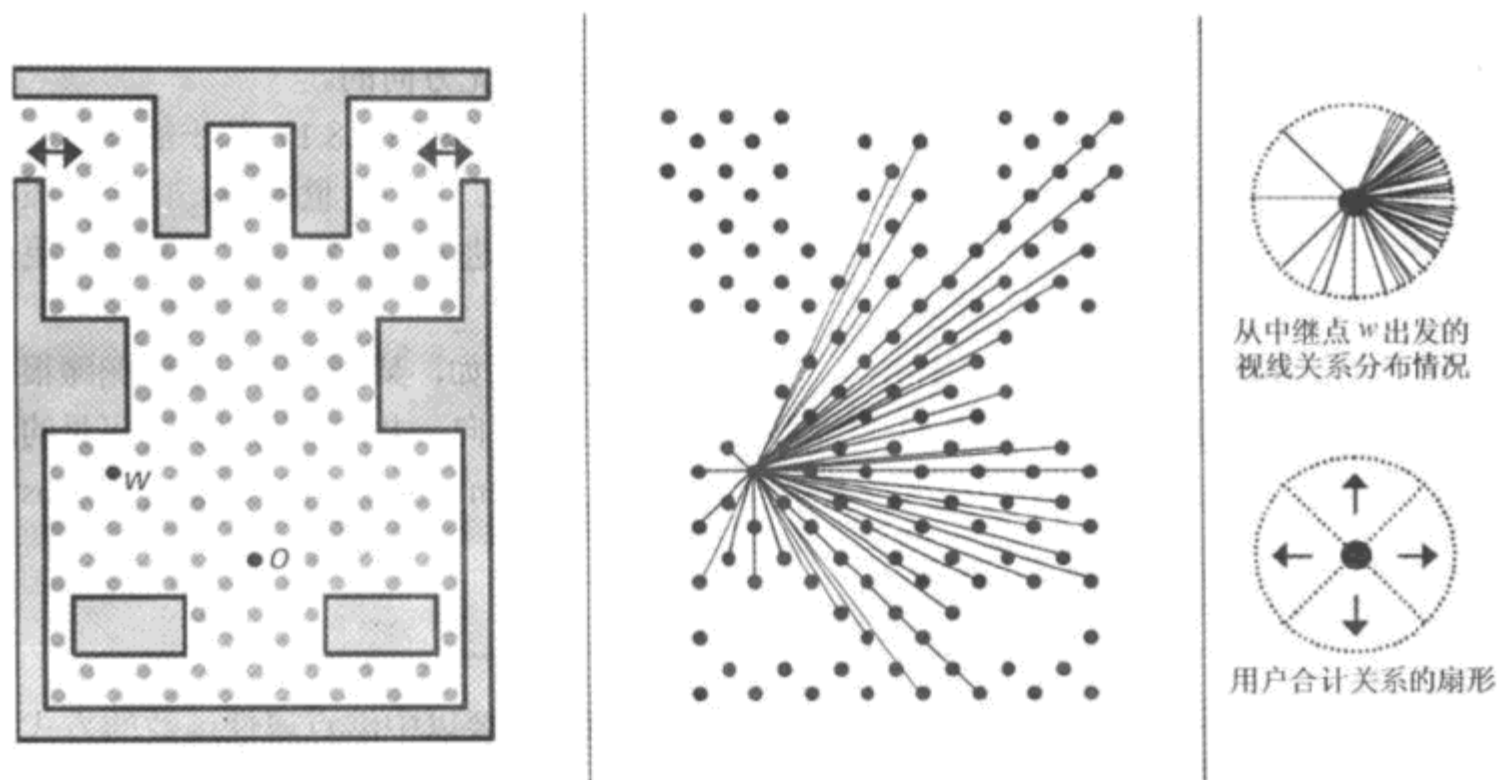


图 3.8.1 (左) 范例地形的俯视图，其中有两个入口/出口和一个目标 (O) 以及大量的中继点 (包括中继点 W)；(中) 范例地形中的中继点以及从中继点到其他中继点的有效视线；(右) 从 W 出发的视线的分布情况以及用于合计关系的扇形

在以非致命武器、缓慢移动、掩体、秘密行动和拦截 (sniping) 为特征的战术模拟中，其他特征至关重要。

- 掩体附近对攻击和防守都至关重要。即使是部分地隐蔽都将是重要的优势。
- 移动缓慢或移动方向可被预测的位置 (如入口) 不是好的进攻点；而俯视这些地方的位置则非常适合防守。
- 与夺旗游戏中一样，俯视目标以及通往目标的主要路径的位置也非常重要。

### 3.8.5 从战术价值到中继点属性

确定大量决定位置的进攻和防守价值的战术特征后，需要将其转换为评估函数和输出。

首先，我们来看看可用于表示战术特征的中继点属性。图 3.8.2 说明了各种可用的中继点属性。

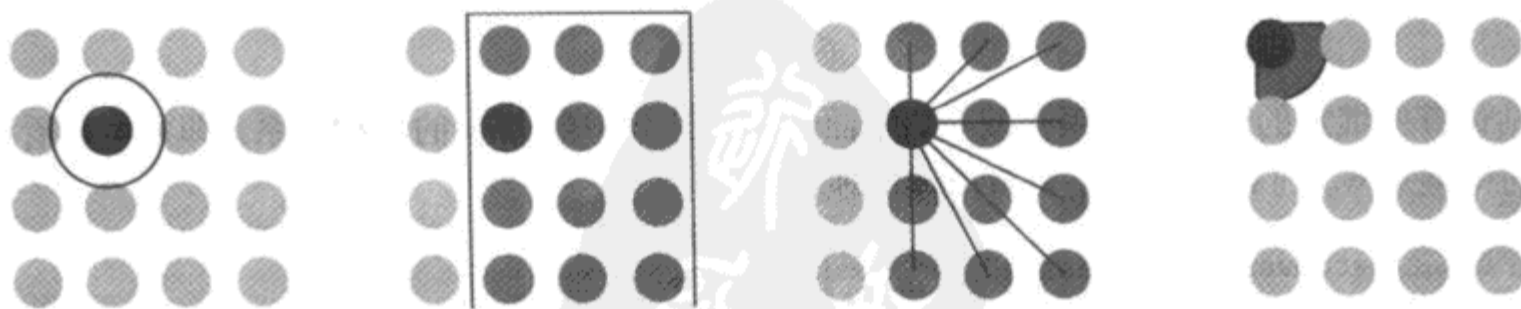


图 3.8.2 中继点属性：(从左到右) 本地属性、组成员、与其他中继点的关系、聚集度

中继点 (waypoint) 包含本地 (local) 属性，如亮度、是否有 power up 或门。另一类属

性取决于中继点是否是大型地形表示（通常是一组中继点）的成员。例如，中继点可能是房间、街道或屋顶的一部分。注意，本地属性和组成员属性都是无方向的。

然而，中继点之间的关系是有方向的。例如，中继点  $w$ （参见图 3.8.1）几乎能够看到其东面的所有中继点，因此从东面接近中继点  $w$  时肯定会被发现。在 3D 世界中，由于海拔高度的不同，从一个方向进入某个中继点很容易，但从其他方向进入将难得多（即需要走更长的路，花更多的时间）。

最后一个（但至关重要）方面是中继点关系的分布情况。例如，某个中继点能够俯瞰很多位于某个方向的其他位置，则玩家或 AI 将能够重点关注这个方向，并立刻发现所有可见的活动，而无需担心来自其他方向的攻击。扇形中的关系密度（concentration）称为聚集度（focus）。

### 3.8.6 计算中继点属性

要为给定中继点和方向计算一个有用的进攻或防御值，需要以中继点属性函数的方式实现每个适用的战术特征。

很多战术特征可被分解为关于中继点图的基本函数，以处理线性距离、行走时间、视线、发射线、目标和障碍物等。例如，好的攻击位置的一个重要特征是，在这里能够快速移动。“水流”位置不允许快速移动，因此本地属性为“水流”的中继点在允许快速移动方面的等级应比较低。

在小房间和隧道中，难以避开火箭和避免火箭或手榴弹爆炸带来的损害，而更开阔的位置在这方面有优势。如果中继点是区域（表示为一组中继点）的一部分，而该区域是一个小房间或隧道，则其在“开阔（open）”方面的等级应较低。因此，中继点的组成员资格可用来计算战术特征。

下面的代码分别使用中继点的本地属性和组成员属性计算了两个战术特征：

```
float GetLocalRapidMovement( waypointid wp )
{ // result in [0 .. 1], higher values meaning higher speeds
  return ( GetActorMovementSpeedAtWaypoint( wp )
           / GetMaxActorMovementSpeed() );
}

float GetOpenAreaMembership ( waypointid wp )
{ // result in [0 .. 1], higher values meaning more open
  return 1.0 - max( IsPartOfSmallRoom( wp ),
                   IsPartOfTunnel( wp ) );
}
```

计算中继点的方向性关系更复杂。下面的函数计算中继点  $w$  的一种关系，即从可以看到  $w$  的中继点  $w_{i0}$  到  $w$  时附近是否有掩体的地方。

```
float GetCoverFromThreatsInDirection(waypointid w) {
  float      property [kMaxDirections];
  unsigned int entry[kMaxDirections];

  // set all property and entry values to 0
  ...
}
```

```

// pass one: collect and interpret relations
for (waypointid w_to = 0; w_to < kMaxWaypointId; w_to++) {
    direction dir = GetDirectionForWaypoints( w, w_to );

    // check for line-of-fire from w_to to w
    if ( (w_to != w) && (HasLineOfFire(w_to, w)) ) {
        entry[dir]++;

        // get value for relation (value in [0..1])
        float value = GetCoverFromThreatsAt( w, w_to );
        property [dir] += value;
    }
}

// pass two: level result into [0 .. 1] range
for ( direction dir = 0; dir < kMaxSectors; dir++ ) {
    if ( entry [dir] > 0 ) {
        property [dir] /= entry [dir];
    }
}

float GetCoverFromThreatsAt(waypointid w, waypointid w_at) {
    for (waypointid w_n = 0; w_n < kMaxWaypointId; w_n++) {
        // check for lack of line-of-fire to neighbor of w
        if ( IsNeighborOf(w, w_n) && (!HasLineOfFire(w_at, w_n)) )
            return 1.0;
    }
    return 0; // no cover found
}

```

第一次遍历所有中继点时，只考虑到  $w$  有发射线的中继点。对于每个这样的中继点， $w$  的关系数增加 1，并将表示附近是否有掩体的值累积起来。

在第二次迭代中，将附近掩体的数量除以关系数，以返回一个 0~1 之间的值。这里没有使用关系数，而是使用函数 `focus()` 处理特定方向的关系集中度。函数 `focus()` 将在后面解释。

在 `GetCoverFromThreatsAt()` 这样的函数（位于上述代码清单末尾的第二个函数）中，进行了非常简单的近似。更高级的近似可能使用模糊逻辑来处理诸如行走到最近掩体所需的时间、可用的掩体数以及武器效力随中继点间距离的变化情况等问题。

中继点的聚集度反映了其不同方向的关系密度。当生存至关重要时，聚集度对于防守尤其重要。在这种情况下，位置仅在某个方向上易于拦截还不够，还需要保护侧翼和后方，也就是说，在两侧和后方的视线/发射线关系较少。

函数 `focus(w, d)` 指出了关系在一个方向相对于其他方向的分布情况。`focus()` 的具体实现取决于需要考虑的方向数（扇形数）以及需要考虑其密度的关系的类型（通常是视线）。图 3.8.3 是范例地形中中继点  $w$  的聚集度实现，它使用 4 个扇形来表示 3D 球形中各种不同的方向。计算聚集度时，应处理这样一种特殊情况，即中继点只在一个扇形中有关系。在这种情



况下，应使用默认的最大值。

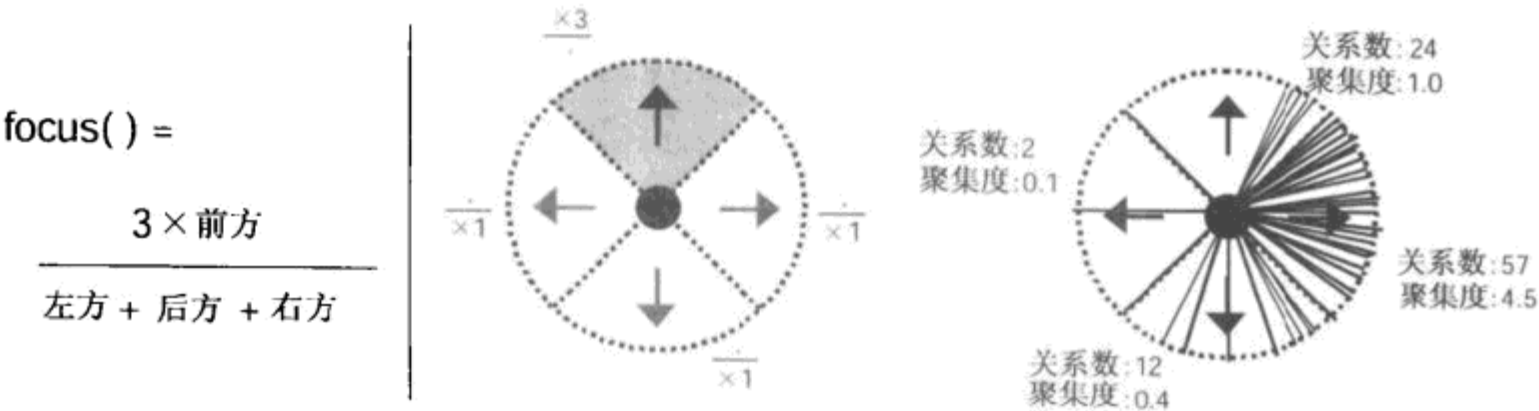


图 3.8.3 以简单的方式计算范例地形中继点  $w$  的聚集度。所有的关系都被投射到一个有 4 个  $90^\circ$  的扇形组成的碟子中。扇形的聚集度是其关系数与其他扇形的关系数的加权比

$\text{focus}()$ 假设中继点在地形中的分布比较均匀，但只要偏离均匀分布不大，该函数也还是健壮的。使用中继点属性计算公式表示所有的战术特征后，便可以以下面的方式将它们合并：

$$\begin{aligned} \text{rating}(w, d) = & \sum k_i \times \text{local\_property}_i(w) \\ & + \sum k_j \times \text{group\_membership}_j(w) \\ & + \text{focus}(w, d) \times \sum k_l \times \text{relation}_l(w, d) \end{aligned}$$

注意，要让  $\text{focus}()$ 正确地强调或抑制  $\text{relation}_x()$ 函数的结果，这些  $\text{relation}_x()$ 函数必须都返回正值。该标称值（ $\text{rating}$ ）指出了 AI 对位置及其在游戏中的作用的先验（ $\text{priori}$ ）理解。该标称值（位置的战术价值）是基于中继点图世界几何形状的。并隐式地使用了一些静态的游戏常量，如行动者的移动速度、武器的性能和  $\text{power up}$  的位置。对于前面讨论的范例区域，我们现在能够自动注释每个中继点在给定方向的进攻或防御值（经过一段时间的试验和调整后）。

### 3.8.7 从经验中学习

显然，并非每个值都完全正确的。通过中继点获得的地形取样就是近似的，用于计算这些值的评估函数也是如此。

另外，我们还忽略了实际的游戏过程。然而，这种错误很容易转化为优势。

由于 AI 使用中继点来执行其行动，并认为玩家位于附近的中继点，因此我们可以很容易地记录玩家在中继点的活动。我们可以以两种方式来使用这些信息。我们可以使用它来校正计算结果以及将其作为计算的额外输入。

例如，可以增加行动者在该中继点的破坏作用和减少其遭受的损害，以提高该中继点的防御价值。换句话说，我们稍微加强 AI 对位置的认识，这样可以校正结果中的部分误差。更重要的是，这样 AI 将能够根据其以往的成功和失败的经验，改变对位置的选择。

捕获的游戏数据还可用作中继点属性的输入。例如，从一个位置看到的敌人频繁通过的中继点越多，该位置的防守价值越大。将游戏数据用作计算的输入时，AI 实际上加深了对地形的认识。



### 3.8.8 将地形推理加入到游戏中

我们的 AI 能够使用几何形状数据、行走时间、最短路径、视线、发射线以及与中继点相关的游戏数据，通过一系列的计算来分析范例区域。这些计算的时间复杂度为  $O(n^3)$ ，因为计算一些中继点到中继点关系时，也需要考虑附近的其它中继点。实际上，这种计算耗时几十秒。这种中继点推理最好作为一个后台线程在预处理关卡（可能在任务之间）完成。

地形推理（terrain reasoning）算法的结果使用的资源很少。通常，每个中继点有多个表，其中每个表包含少量的字节。这些表可被快速地读取，它们包含有关地形的知识，而这些知识本来几乎是无法获得的。另外，基于游戏数据获得的认识使变化多端和自适应的 AI 行为得以实现，而其成本是微不足道的。

通常，基于中继点的推理并非是 CPU 密集型的。只要 AI 考虑的中继点为数不多且有中间推理结果，地形推理在游戏期间是切实可行的。对于每个中继点，使用小型查找表来存储附近的中继点和可见的中继点，AI 能够高效地规划路径，避开位于敌人火力范围内的位置。

动态游戏地形和实体，如门、车辆、可被破坏的地形以及烟雾，使地形推理更复杂，因为这些动态变化使得部分预先计算的结果无效。尤其是视线、发射线和路径，它们将受到动态地形的影响。尽管如此，但使用预先计算的中间结果还是很有吸引力的。与根本不进行推理相比，使用这些结果并针对动态变化进行校正，智能化程度和效率将更高。

这里讨论的地形推理是一组直观推断。为 AI 地形推理找到合适的因素和权重，需要进行大量的实验并进行一些分析。为高效地调整和验证您的算法，必须以可视化的方式表示推理的中间结果。将结果投影到游戏的地形中，将更容易了解它们。也可以将结果导出到电子表格中，以便做进一步的分析。

### 3.8.9 其他应用

这里使用一个简单的范例，预先计算游戏位置的进攻和防御价值，并介绍基于中继点的推理。推理概念的应用范围要比这个范例宽得多。剖析每个位置的战术方针并将其表示为评估函数这一方法可用于很多 AI 需要对位置进行推理的游戏。

请自问下面的问题：

如果要保卫这个基地（或 AI 需要做的任何工作），为何选择位置  $x$  而不是位置  $y$ ？距离、行走时间、区域类型（或任何地形属性和您能想到的任何游戏物体）有何影响？

被放在密度足够高的、跨越地形的图形中时，中继点不仅仅可用于导航，还有很多其他的用途。结合使用针对附近中继点和可见中继点的小型查找表时，中继点让 AI 能够预测看不见的对手的移动情况，进而在附近找到一个能够看见对手的位置。

可通过中继点在 A\*寻径算法中加上额外的成本。如果将从设想的威慑位置能够见到的所有中继点加上额外的成本，则 A\*寻径程序将提供一条可以避开威慑的路径。如果同时对限制移动的中继点进行标记，则返回路径的战术性将更强。这只是众多基于中继点推理的应用范例中的两个。您可能还会想到其他一些使 AI 处于优势的例子。

### 3.8.10 结论

---

中继点图提供了一种易于使用的游戏地形表示。可以将游戏中很多与地形相关的战术转换为与中继点相关的属性。例如，本文介绍了如何创建一个中继点评估函数，以表示位置的进攻价值、防御价值以及对诸如埋伏和拦截等特定战术的价值。这将让 AI 对地形有深刻的认识。

可以使用捕获的游戏数据来扩展和改进评估函数，让 AI 更为变化多端，自适应性更强。AI 还可以采用很多其他方式，使用中继点对地形进行推理。只需将战术方针同（存储在查找表中的）中继点属性关联起来即可。

### 3.8.11 参考文献和推荐读物

---

[Pottinger00] Pottinger, Dave, "Terrain Analysis in Realtime Strategy Games," Proceedings of Computer Game Developer Conference, 2000.

[Rabin00] Rabin, Steve, "A\* Speed Optimizations," *Game Programming Gems*, Charles River Media, 2000: pp. 272~287.

[Reece00] Reece, Doug, et al, "Tactical Movement Planning for Individual Combatants," Proceedings of the 9th Conference on Computer Generated Forces and Behavioral Representation, 2000. Also available online at [www.sisostds.org/cgf-br/9th/](http://www.sisostds.org/cgf-br/9th/).

[Snook00] Snook, Greg, "Simplified 3D Movement and Pathfinding Using Navigation Meshes," *Game Programming Gems*, Charles River Media, 2000: pp. 288~304.



## 3.9 用于可视点寻径的扩展几何体

Thomas Young

Thomas.young@bigfoot.com

在《游戏编程精粹 1》中, Bryan Stout 和 Steve Rabin 介绍了可视点(points-of-visibility)寻径(pathfinding) ([Stout00]、[Rabin00])。这是一种找出避开多边形障碍物的最短路径的方式。

正如 Steve 指出的, 这种方法有很多重要的优点。它让您能够创建搜索空间的最小表示, 因此搜索速度非常快。另外, 找出的路径非常直。

本文将介绍如何将根据多边形障碍物构建扩展几何体的工作自动化进程以及如何使用扩展集几何体实现可视点寻径 (Points-of-visibility pathfinding)。通过使用扩展几何体克服了这种技术的缺点。

- 可直接从扩展几何体中提取可视点, 而无需设计器的帮助;
- 支持动态物体 (如角色、门、移动的积木等), 这是通过动态地生成和集成这些物体的扩展几何体实现的;
- 可参数化扩展以支持不同大小和构造的角色。

在很多游戏中, 寻径对 AI 来说是必不可少的。角色行为的方式是由寻径程序(pathfinder)保证的。例如, AI 角色不能因障碍物而无法前进。角色必须能够找到绕开障碍物的路径 (如果该路径对玩家而言是显而易见的话)。当与障碍物的碰撞非常复杂时, 即使是这些能力也难以保证。

如果我们接受角色移动的碰撞模型的一些局限性, 则可以将相同的模型用于碰撞和寻径。通过使用扩展几何体, 可以创建一个对这种碰撞模型进行精确校正的寻径系统。该系统能够正确地理解角色能够到达的任何位置, 并返回一条畅通无阻的路径, 从而确保角色不会因障碍物而无法前进。

基于“完美”寻径系统的高级 AI 的代码更容易编写, 因为您无需检测那些微妙的情形, 如角色因为无法找到从当前位置出发的路径而受阻。通过完全依赖于寻径系统的结果, 可以创建复杂的、与移动相关的 AI, 而需要的代码行少得多。这也更有趣。

### 3.9.1 定义碰撞模型

碰撞模型精确地定义了角色同环境的碰撞。我们基于每个角色的“碰撞形状 (collision shape)”来创建碰撞模型。碰撞形状是一个能够最好地表示角色的大小和形状的凸多边形。当角色移动时, 该形状将用角色的原点来表示, 但当角色转身时, 该形状不会旋转。在角色的碰撞形状同环境

重叠的位置，角色将不能前进。

环境用一系列的多边形表示，其中包括表示受阻区域的凸多边形和非凸多边形、表示畅通区域的 2D 网格（mesh）或这些形状的组合（见图 3.9.1）。

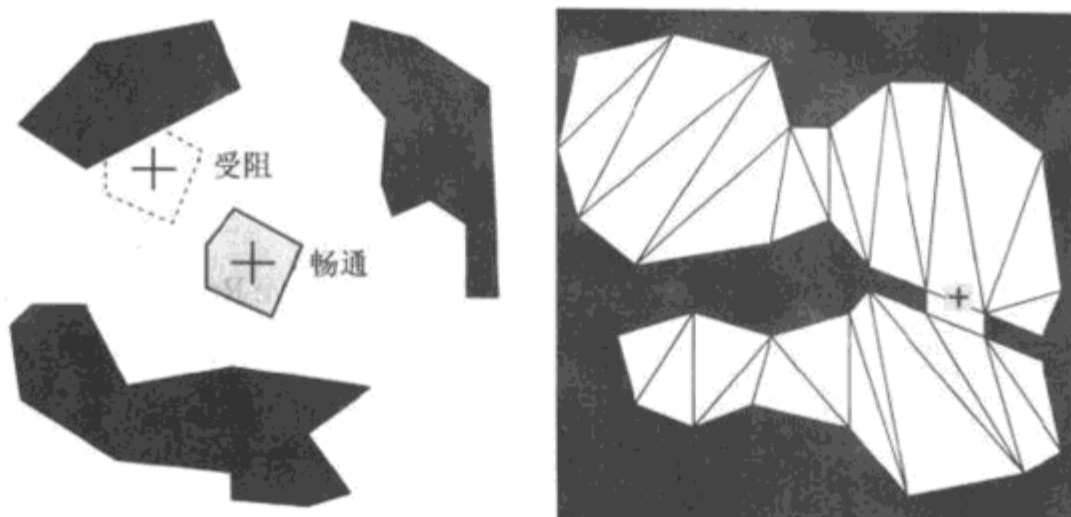


图 3.9.1 多边形环境中的碰撞形状

### 3.9.2 多边形寻径

定义碰撞模型后，接下来由寻径程序为角色找出一条在该模型中畅通无阻的路径。出于简化的目的，我们只考虑传统的寻径约束条件，即找出从起点位置到目标位置的最短路径。可以用另一种方式来描述这个问题，即找出一组畅通无阻的、将起点和目标连接起来的线段。

### 3.9.3 扩展并解决问题

这里的诀窍是，创建一个将碰撞形状（collision shape）和环境中的障碍物形状组合起来的扩展几何体。这种表示极大地简化了寻径系统所需的查询。

我们使用的扩展是平面集的闵可夫斯基和（Minkowski sum）。具体地说，扩展几何体是环境与非（negated）碰撞形状的和，这有时候被称为闵可夫斯基差（Minkowski difference）。

集合 A 和 B 的闵可夫斯基和是一组这样的值，即不能将集合 A 的某个成员和集合 B 的某个成员相加得到。多边形可视为一个平面集，即由多边形内的点组成的集合。

我们的扩展几何体表示这样一些点组成的集合，即可以通过将环境中的某个点减去碰撞形状中的某个偏移量（offset）得到。这意味着对扩展几何体中的每个点，碰撞形状中都有一个与环境重叠的点与之相对应。因此，扩展几何体表示一组这样的点，即在这一点，角色将受到阻碍（如图 3.9.2 所示）。

碰撞形状同多边形环境之间的碰撞与碰撞形状中的某个点同扩展几何体之间的碰撞等效。为判断角色是否能够进入某个位置，我们检测该位置是否位于扩展几何体中；为判断角色是否能够沿某条线段前进，只需检测该线段是否同扩展几何体发生碰撞。

对于可视点寻径（points-of-visibility pathfinding），扩展几何体中的一系列凸角（convex corner）提供了可视点。如果两点之间的线段不会同扩展体碰撞，则这两点之间便是可以连接的。

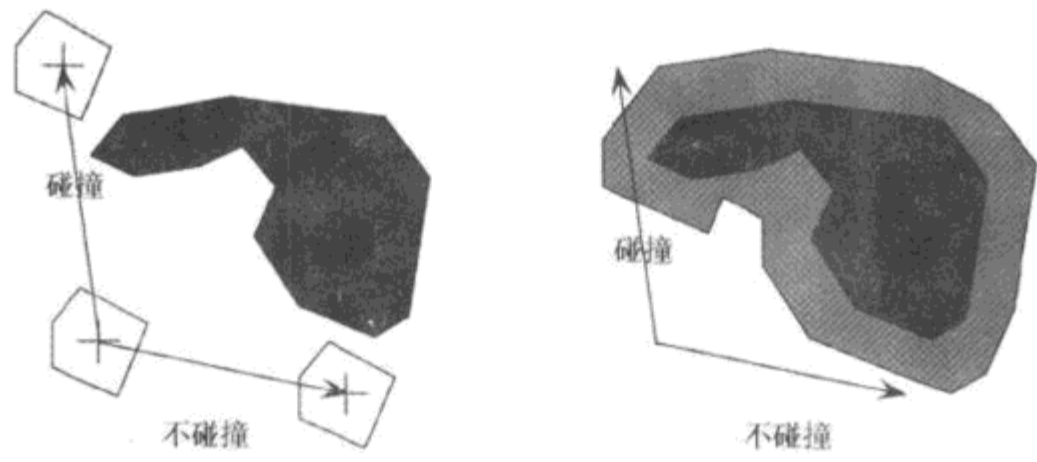


图 3.9.2 多边形环境中的碰撞形状和扩展几何体中的点

3.9.4 凸多边形的闵可夫斯基和

为创建扩展几何体，首先来看一种比较简单的情况，即障碍物是由单个凸多边形组成的。碰撞形状的非也是一个凸多边形，因此我们必须计算两个图多边形的闵可夫斯基和。

对于凸多边形 C 和 O，可将闵可夫斯基和视为让 C 的中心沿 O 的周边线移动所扫过的空间（如图 3.9.3A 所示）。该空间位于一个更大的多边形内。也可将闵可夫斯基和视为将碰撞形状沿 O 的周边线移动时，其中点的轨迹组成的性状（如图 3.9.3B 所示）。

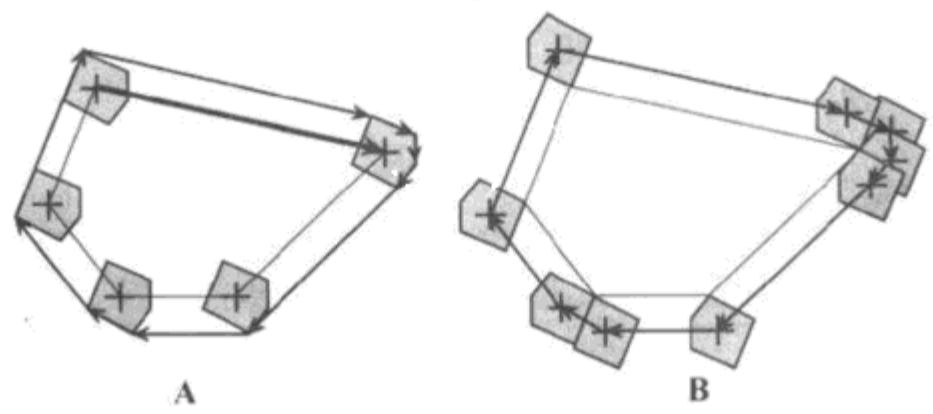


图 3.9.3 两种看待凸多边形的闵可夫斯基和的方式

扩展多边形的边通过下列三种方式来生成（参见图 3.9.3B）：

- (1) 当 C 位于 O 的顶点处时，直接来自 C 的一条边；
- (2) 当 C 沿 O 的边前进时，由 C 的一个顶点生成；
- (3) 当 C 沿 O 的平行边前进时，由 C 的一条边生成。

如果没有平行边，则 C 和 O 的每条边都只被使用一次，这样扩展几何体的边数将等于 C 和 O 的边数和。这对于预先为扩展几何体分配存储空间很有用。

创建扩展几何体非常简单。按顺时针方向对 C 的顶点进行编号（如图 3.9.4A 所示）。对于 O 中的每条边，判断使用哪个顶点来扩展这条边的起点和终点（如图 3.9.4B 所示）。这样得到的将是(2)类和(3)类边。当某条边的终点不与下一条边的起点相连时得到的是(1)类边，必要时对 C 中的顶点进行插值，并添加边（参见图 3.9.4B）。插值只不过是依次考虑 C 中的顶点，直到找到扩展下一条边的起点的顶点为止。

那么，开始时如何确定每条边的起点扩展和终点扩展呢？可以根据 C 中的边的方向，定



义一个基于向量的循环序列（参见图 3.9.5）。对于  $O$  的每条边，该序列的位置指出了应使用  $C$  的哪个顶点来扩展这条边的起点和终点（参见图 3.9.5）。为提高速度，可以在遍历  $O$  的边时依次记录该位置。

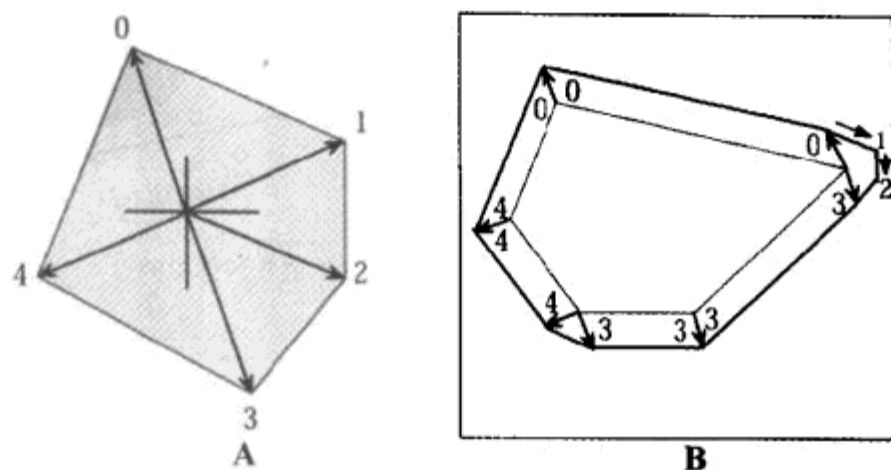


图 3.9.4 A)  $C$  的顶点；B) 起点和终点（插入的顶点）

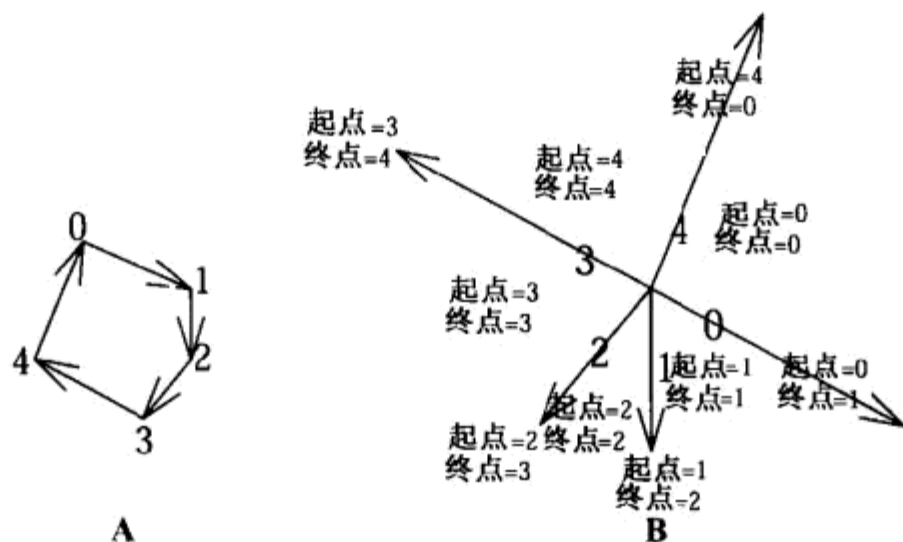


图 3.9.5 按边向量排序以及边的起点扩展和终点扩展

### 3.9.5 扩展非凸几何体

能够扩展凸多边形后，便可以很容易地将这种方法扩展到非凸多边形。只需将非凸多边形分割成多个凸多边形，然后分别对每个凸多边形进行扩展即可（参见 3.9.6A）。然而，这将得到大量相互重叠的多边形，其中的很多边和顶点是多余的。

为获得格式良好的闵可夫斯基和，必须检测扩展多边形之间的交点，并将这些多边形连接成一个扩展形状。这种方法存在的一个主要问题是，交点可能无法用我们当前用来表示坐标的数字来表示。另外，如果用近似的值来表示交点，则我们的寻径程序将不能与碰撞模型协调一致，虽然如果根据同样的扩展表示来创建碰撞引擎则这将无关紧要。

一种不错的解决方案是计算“懒人（lazy man）闵可夫斯基和”。这里我们关心的不是阻碍点集合，而是从畅通区域到非畅通区域的过渡。对于外边界无法被分解为凸多边形的情况，这种方法更合适。另外，如果我们以后需要进行扩展，以支持诸如重叠几何体等特性时，这



种方法的关联性也更强。

要创建“懒人闵可夫斯基和”，可以使用几乎与扩展凸多边形相同的方法。差别在于，在凹角处不插入新的点。如果凹角前面的边的终点与其后面的边的起点不相同，则在这个地方，扩展后的边将不是首尾相连的，且会有一个交点（如图 3.9.6B 所示）。

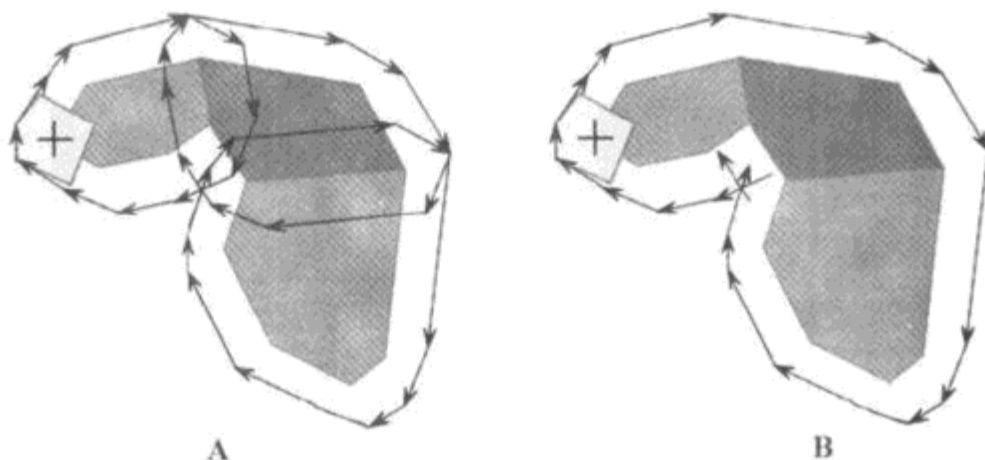


图 3.9.6 A) 将非凸多边形划分为多个凸多边形并分别进行扩展；B) 懒人闵可夫斯基和

### 3.9.6 选择碰撞形状

碰撞形状不能旋转是我们碰撞模型的一个主要缺点。由于游戏中的角色几乎总是需要旋转，因此我们希望使用可旋转的对称碰撞形状，以实现最为一致的碰撞行为。圆是最佳的近似，但我们不能使用圆，因为这样扩展几何体将没有角，因此也没有可视点。我们将使用  $N$  边正多边形来逼近圆。

通常，碰撞形状的边越多，扩展几何体的边和角也越多。边角也越多，寻径的成本将越高（有时呈几何级数增长）。正方形和八边形是显而易见的选择，因为它们相对简单且可以与坐标轴平行。

以前，为提高性能，碰撞形状与坐标轴平行至关重要。这样将向量同水平线、垂直线或  $45^\circ$  角的线进行比较时，常常无需进行乘法运算。现在，分支预报问题意味着不论在何种情况下都执行同样的乘法运算更合适，因此不再因为性能方面的原因而只使用与坐标轴平行的形状。然而，对于有些平台，避免乘法运算仍然是很重要。

### 3.9.7 结论

本文介绍了如何构建一个这样的寻径系统，即对一个合理而有趣的多边形碰撞模型而言它非常精确。必须在复杂的碰撞系统和角色能够理解的碰撞系统之间进行折衷。如果能够在游戏中将这种模型用于角色碰撞和寻径，您将获得巨大的回报——减少调试工作并简化 AI。

### 3.9.8 参考文献

[Rabin00] Rabin, Steve, "A\* Speed Optimizations," *Game Programming Gems*, Charles River Media, 2000.

[Stout00] Stout, Bryan, "The Basics of A\* for Path Planning," *Game Programming Gems*, Charles River Media, 2000.



## 3.10 优化可视点寻径

Thomas Young

thomas.young@bigfoot.com

“可视点”寻径 (pathfinding) 方法有很多优点 (参见《游戏编程精粹 1》中 Steve Rabin 撰写的文章, [Rabin00])。另外, 对于给定的多边形障碍物集合, 通过使用扩展几何体可使该方法获得正确的路径 (参见本书中“用于可视点多边形寻径的扩展几何体”一文)。其复杂度与规模 (scale) 无关, 因此能够高效地计算跨越很长距离的路径; 同时, 通过预先计算可视点之间的连接网络, 可以构建速度非常快的寻径程序。

随着地图越来越复杂, 可视点的数量也越来越大, 而可视点之间的连接数将呈几何级数增长, 尤其是当地图中包含大型的空旷区域 (open area) 时。内存将成问题, 尤其是在控制台平台中。另外, 由于要求 AI 有趣, 因此寻径程序必须支持动态特性, 如移动的物体、角色、门等。为支持这样的动态特性, 必须使用动态生成的可视点以及能够移动进而导致连接无效的障碍物。如果不小心, 在支持这些动态特性的情况下计算可视性的开销也将呈几何级数增长。

本文介绍一些优化措施, 让我们能够迅速放弃大量潜在的路径, 极大的缩小连接网络的规模。有了这些优化措施, 我们的算法在处理日益复杂的地图时可扩展性将更佳。

### 3.10.1 可视点寻径 (points-of-visibility pathfinding)

本文假设寻径主体可被视为多边形环境中的一个点。前一篇文章介绍了如何将一种更有趣的碰撞模型转换为这种形式。

我们的寻径程序使用 A\* 算法, 并将一组可视点作为从起点到目标之间的可能状态。可视点是直接从多边形环境中的凸角派生而来的。有关 A\* 算法的介绍, 请参阅《游戏编程精粹 1》中 Bryan Stout 撰写文章 [Scout00]。

### 3.10.2 存储到每个点的最短路径

为防止搜索树呈几何级数增长, 首先需要确保只保留到中间点的最短路径, 供以后进一步的判断。这并非可视点寻径特有的优化, 因为这可能是任何 A\* 实现的标准行为。然而, 这里有必要再次强调它, 因为它将给可伸缩性带来翻天覆地的变化。

Bryan Stout 在其文章中描述了如何使用开链表 (open list) 和闭链表 (closed list)。在 A\* 算法的每一步中, 都需要搜索这些链表, 找出其他所有以当前考虑的点为终点的路径。

基于分片 (tile) 的寻径和可视点寻径之间的一个重要区别是, 需要考虑的可能的中间位置的数量。由于我们只考虑环境中的凸角, 因此其数量级常常小于基于分片的寻径 (tile-based pathfinding) 的数量级——表示地图所需的分片数。

这意味着可以使用一个数组 (每个元素对应一个可视点) 来记录到目前为止找到的到每个点的最短路径。我们无需搜索开链表和闭链表, 而只需在数组中进行一次查找。事实上, 这意味着根本不需要开链表和闭链表。

每次使用该数组前必须将其清空, 但在大多数平台上都有现成的内存清空函数, 因此这种操作所需的时间通常是无关紧要的。如果有大量的点, 而您希望每次搜索只使用其中的很小一部分, 则可能有必要维护一个闭链表, 以记录搜索结束后必须清空哪些点。

### 3.10.3 将凸角相连

在 A\* 的每一步中, 都需要基于目前找到的最佳部分路径生成一组后续路径。每条后续路径都是通过使用一条线段将部分路径同另一个点相连得到的。这个点可能是可视点, 也可能是目标点。直观的方法是, 对于对应的线段不会被阻断的每个点, 都生成一条后续路径。

通过预先处理任何一对可视点之间的线段同环境的碰撞情况并创建一个表 (可能很大), 可以在运行阶段快速地判断出从某个源可视点能够达到其他哪些可视点。在运行阶段, 起始位置和目标位置将发生变化, 因此无法预先处理这些点之间的线段。对于动态障碍物导致的可视点, 情况也是如此。

#### 1. 优化 1: 只考虑到轮廓点

正如从某些源点看到的, 每个点都可被归类为做轮廓 (silhouette) 点、右轮廓点或非轮廓点 (如图 3.10.1 所示)。

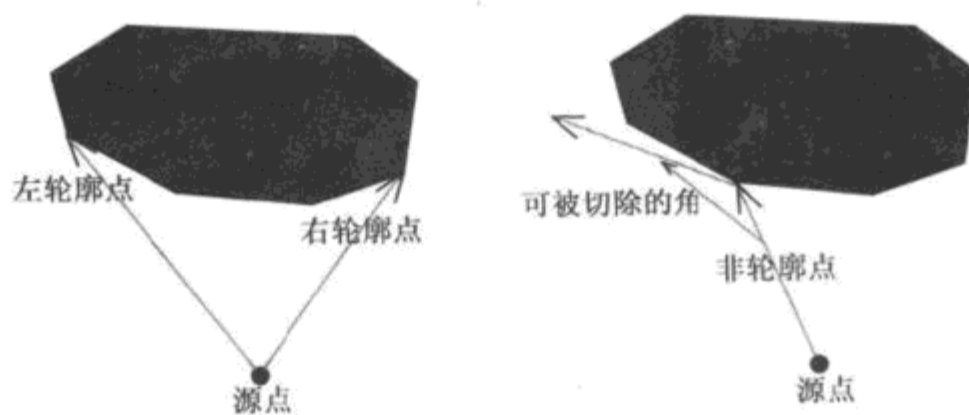


图 3.10.1 轮廓点

这种优化是基于观察结果的, 即无需考虑连接到非轮廓点的路径段。任何到非轮廓点的路径段都将导致一个拐角, 可以切除该拐角获得更短的路径。只要该点不与目标点或另一个轮廓点重叠, 上述结论都是正确的。对于这两种情况, 可以采取另一种机制来生成路径。

这种优化适用于生成部分路径的一组后续路径。这里将部分路径的终点称为当前点; 我

们考虑将其用于扩展路径的点是潜在的下一个点。我们将丢弃所有这样的点，即从当前点看它们不是轮廓点。

我们之所以丢弃非轮廓点，这是因为对于使用这样的点得到的路径，总是可以通过切除一个拐角使之更短。为更好地了解这里发生的情况，知道这样一点将有所帮助，即随着搜索的进行，终将找到这条更短的路径。这发生在将当前点直接连接到目标点或连接到障碍物上的轮廓点时。

如果使用可视性图 (visibility graph)，这种优化也可直接用于其中的连接。任何从源点看，终点不是轮廓点的连接都可将其从图中删除。

## 2. 优化 2：只考虑绕角前进的路径

这种优化也适用于生成部分路径的后续路径。

在这种优化中，我们关心的是部分路径中的最后一条线段。从起始位置开始，首次生成后续路径时，部分路径为空，因此这种优化不适用。这里将部分路径中最后一条线段的起点称为前一个点；将部分路径的终点称为当前点。作为优化 1 的结果，从前一个点看，当前点总是轮廓点。

这种优化后面的推理类似于第一种优化后面的推理。任何不绕轮廓点前进的路径都将导致一个拐角，通过切除该拐角可以得到一条更短的路径。图 3.10.2 显示了一个轮廓点、一组绕该轮廓点前进的路径段以及一条可以丢弃的路径。同样，最短的路径要么是直接到目标点的连接，要么将经由另一个轮廓点（如果直接连接被阻断）。

为实现这种优化，我们使用从前一个点到当前点的向量以及从当前点到当前考虑的潜在的下一个点的向量。在左轮廓点处，下一个向量必须位于前一个向量的右边，且在障碍物的左边。在右轮廓点处，下一个向量必须位于前一个向量的左边，且位于障碍物的右边（如图 3.10.2 所示）。

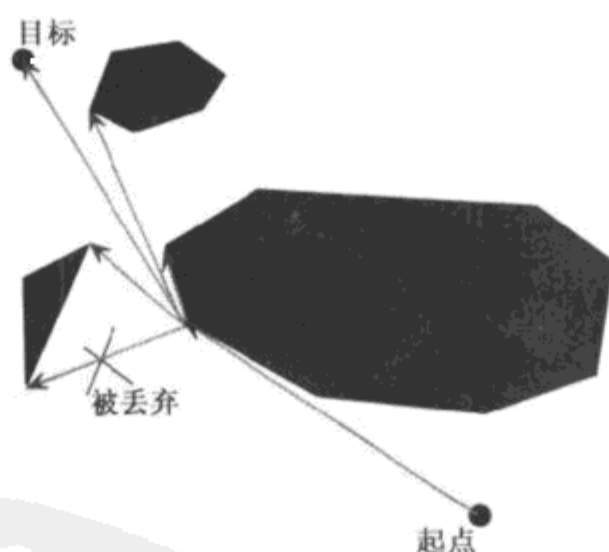


图 3.10.2 路径必须绕每个角前进

### 3.10.4 轮廓区

现实时，对于每个可视点定义两个与之相关轮廓区 (silhouette zone) 很有帮助 (图 3.10.3)。每个可视点都位于凸障碍物的角上，而两个轮廓区是通过延长该角前后的障碍物边得到的。

从某个点看时，轮廓区将可视点归类。如果该点位于左轮廓区，则该可视点为左轮廓点；如果该点位于右轮廓区，则该可视点为右轮廓点。

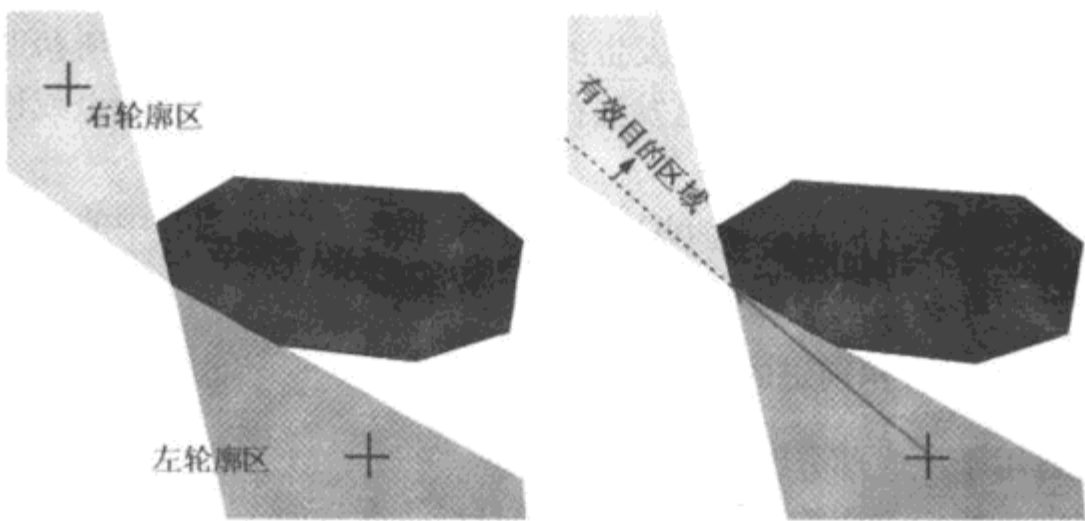


图 3.10.3 顶角处的轮廓区

为绕轮廓点前进，下一个点必须与前一个点位于相反的轮廓区，且必须位于有效的目的区域。有效目的区域以前一条连接的延长线为边界。

可以将第二种优化用于可视性图中的连接。图中的每条连接都有源点和目的点。生成后续连接时，这两个点分别是当前点和下一个点，但在简单的可视性图中，没有有关前一个点的位置的信息。这意味着将这种优化用于可视性图时，必须考虑前一个点的所有可能位置。这意味着有效目的的区域由两个轮廓区组成。

因此，要将这种优化用于可视性图，还需丢弃其目的点不位于源点的任何一个轮廓区中的连接。当我们检索可视性图中的连接时，知道前一个点的位置，因此可以做更具体的检测，以判断是否应丢弃该连接。

表 3.10.1

	点数	优化前的连接数	优化后的连接数
环境 1	21	231	98
环境 2	96	1638	568

图 3.10.4 是两个障碍物环境的例子。上表是使用轮廓优化丢弃连接之前和之后，可视性图中的连接数。

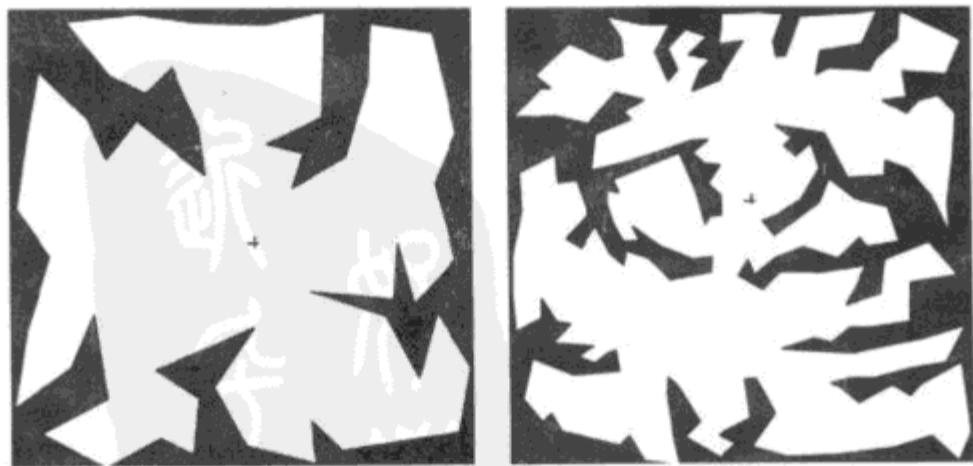


图 3.10.4 环境 1 和环境 2



### 3.10.5 将轮廓区用于空间分区系统

---

本文介绍的优化让我们能够减少可视性图中的连接数，进而缩小搜索树。对于动态点，仍需要检查与其他每个点之间的连接。我们能够迅速丢弃其中的很多连接，但对于不能丢弃的连接，必须对其碰撞情况进行检测。

接下来，需要创建一种表示，以快速确定最小的、可能与给定的动态点相连的点集。对于创建这种表示，轮廓区是个不错的起点。给定可视点的轮廓区提供了一个区域，位于该区域的动态点可能能够连接到该可视点。这种表示的具体特性将随障碍物环境的类型而异。

在可见性较低、由通过门（portal）连接起来的房间组成的地图中，可通过门对轮廓区进行投影，以判断哪个房间可能能够看到相应的可视点。

对于更通用的系统，可根据多边形环境对轮廓区进行裁剪，以确定能看到给定点的区域。然后，将得到的形状输入到某种通用的空间分区系统。如果这些区域的形状被准确地表示，则只要知道动态点位于该形状中，便无需进行碰撞检测（注意，要准确地表示通过裁剪得到的区域，可能需要准确地表示直线的交点）。

### 3.10.6 结论

---

要构建一个高效的可视点寻径系统，涉及的细节比这里介绍的要多得多。本文介绍了一些迅速丢弃很多连接，进而减少需要考虑的路径总数的技术。这是一个良好的开端。

本文还简要地介绍了如何将轮廓区用于空间分区系统，以有效地处理动态点。为有效地支持动态障碍物，必须使用优良的空间分区系统。

最后是检测动态点到其他点的连接问题。对于非动态障碍物，只需检测一次，即从起点到目标点的潜在连接；对于动态障碍物，则可能需要进行大量的这种检测，因此高效的实现就显得至关重要。

### 3.10.7 参考文献

---

[Rabin00] Rabin, Steve, "A\* Speed Optimizations," *Game Programming Gems*, Charles River Media, 2000.

[Stout00] Stout, Bryan, "The Basics of A\* for Path Planning," *Game Programming Gems*, Charles River Media, 2000.



## 3.11 有齿物群的模拟：捕食者和猎物

Steven Woodcock

ferretman@gameai.com

**物**群模拟 (flocking, 有时也叫作 swarming 或 herding) 技术是由 Craig Reynolds 在 1987 年的 SIGGRAPH 会议中发布的论文 *Flocks, Herds and Schools: A Distributed Behavioral Model* [Reynolds87] 中首次提出的。在该论文中, Reynolds 定义了物群模拟 (flocking) 的 3 种基本规则 (或导航行为, steering behavior), 并解释了这些规则如何相互作用, 给名为 boid 的飞禽以逼真的群体行为。在《游戏编程精粹 1》中, 作者撰写了一篇短文, 对这一主题做了简要的介绍, 并添加了另一种导航行为, 作者将这些导航行为称为“物群模拟的 4 种规则”, 它们分别是:

- 分离 (separation): 同物群中的其他成员若即若离;
- 列队 (alignment): 与物群中的其他成员保持相同的航向;
- 内聚 (cohesion): 不掉队;
- 避开 (avoidance): 避开障碍物和天敌。

对于这 4 种简单的行为规则, 令人感兴趣的是, boid 的最终行为有多逼真。从《游戏编程精粹 1》提供的演示程序 (也包含在本书附带的光盘中) 可知, boid 聚集在一起, 并以波浪形飞行。遇到另一群 boid 时, 它们将逃开, 并在必要时分成多群, 以避免与这些 boid 接触。分开后, 如果 boid 找到其他的伙伴, 它们将组成新的 boid 群, 并最终找到并加入其他的 boid 群。

物群模拟技术另一个有趣的方面是, 移动算法本身是无状态的——在移动更新中, 不记录任何信息。在每次更新循环中, 每只 boid 都将重新评估其环境。这不断降低了内存需求 (如果使用其他方法来提供类似的行为, 则需要更多的内存), 同时让物群能够对不断变化的环境状况做出实时的反应。这样, 物群将具备突发行为 (emergent behavior) 的特性, 即物群中的所有成员都对要前往何方一无所知, 但作为一个整体行动, 避开障碍物和天敌, 并动态地保持若即若离。

实践表明, 这种技术对视频游戏很有用, 它用于从 RTS 游戏中的部队编队到 RPG 游戏中的群体逼真行为等各项内容中。

本文以《游戏编程精粹 1》中的那篇文章为基础, 完成了其中建议的一些改进, 并揉和了来自读者的一两个建议。作者将在原有的演示程序中添加一些特性, 使 boid 各不相同, 其中的一个特性是饥饿。同时环境中加入一些障碍物, 使 boid 导航起来更困难些。然后加入一个新的鉴别器,

让 boid 各自逃命——有些 boid 是捕食者，以其他 boid 为食！boid 将分为 3 类：老鹰、麻雀和昆虫。为此，需要加入第五个规则：

- 生存 (survival)：必要时进行捕食或发现捕食者时逃脱被吃的命运。

3.11.1 全新的世界

原有演示程序中的立方体世界乏味而空洞，除了 boid 外一无所有。boid 几乎将不断地漫无目的地游荡（假设它不会变成其他 boid 群的成员）。在与物群模拟相关的代码中，惟一可能对 boid 的飞行有影响的代码是物群模拟算法中的一个较小的函数 `CBoid::Cruising()`，然而其作用不大。约束世界的代码（让 boid 从世界的一边飞到另一边）本可以让物群更有趣些，但物群成员最终总是会再次找到其他成员。

要使这些飞禽的生活更有趣，必须增加种类，并提供一个与实际游戏更接近的环境。为了让 boid 更多地考虑周边环境，我们增加需要处理的新东西——障碍物。

1. 障碍物

顾名思义，障碍物指的是挡道的东西。为创建障碍物，增加了一个名为 `CObstacle` 的类。`CObstacle` 对象形成一个无法穿越的锥形壁垒，boid 在前进时将尽可能地避开它们。

这个类提供了一种让世界更有趣的简单途径，并给 boid 的飞行带来了一些有趣的挑战。可以在世界的任何地方创建障碍物，并任意指定其朝向；创建障碍物时，可以使用指定的值，也可以使用随机生成的值。新增的方法 `CObstacle::GetPosition()` 用于碰撞检测，并判断障碍物是否在给定 boid 的视线范围内；而方法 `CObstacle::GetRelativeDiameter()` 将返回对于 boid 的巡航高度而言，障碍物的相对大小。

2. 各种飞行动物

图 3.11.1 表明，boid 现在分为 3 类：老鹰、麻雀和昆虫，它们都是飞行动物。这些动物分别在本文使用的生态环境中扮演一个重要角色，它们捕食其他动物。

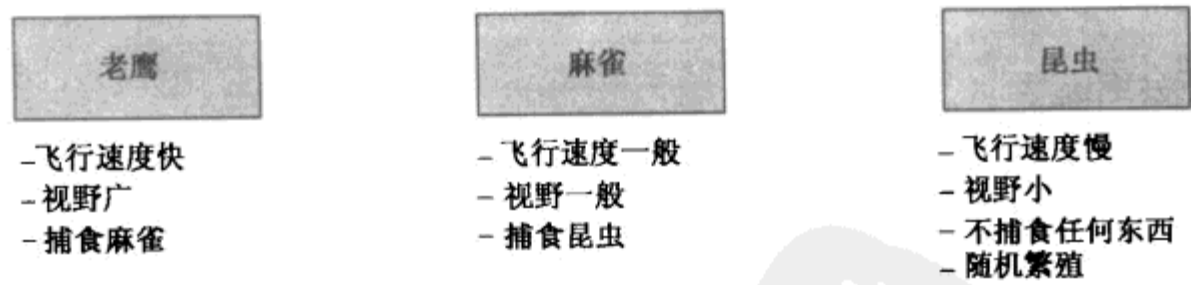


图 3.11.1 boid 的类型

3. boid 各不相同

原有的演示程序 `SimpleFlocking` 初始化 boid 的各种参数——视野、最大飞行速度等，所有 boid 的这些参数值都相同。除了所属的物群外，没有专门的鉴别器将 boid 区分开来。其他物群中的 boid 被视为天敌，应像逃避瘟疫一样避开它们（如果该开关被打开的话）。

演示程序 `PredatorsandPreyFlocking` 在创建 boid 时将随机生成的参数值赋给它们，使 boid

各不相同。这样，（使用 `CBoid` 类中新增的构造函数）创建的 `boid` 将具备某些个性化的东西，这使得每个 `boid` 都与其伙伴有些不同。有些 `boid` 的视野小，而有些比其他 `boid` 的视野广；同样，有些 `boid` 可能希望离同伴比标准距离更远；有些 `boid` 可能比其他 `boid` 更饥饿；等等。

为何要这样做呢？原因有两个，旨在使生物的行为更逼真。首先，相对于每只 `boid` 都相同，每值 `boid` 有不同的能力更接近现实；其次，这些不同使得 `boid` 交互时将呈现突发行为，这也使得集体行动的物群更逼真。同一物群中的两只 `boid`（其中一只希望同伙伴保持较近的距离，而另一只希望更远些）之间忽近忽远将形成有趣的群体动态行为；同样物群中某个 `boid` 比其他伙伴的视野更广，能够更早地发现捕食者，也使群体动态行为更有趣。

另外，新的演示程序可以更精确地控制哪类 `boid` 将被视为天敌，哪些类型不是。在 `CBoid` 的构造函数中，使用了一个新参数来指定 `boid` 的类型。这样，可以在方法 `CBoid::SeeEnemies()` 和 `CBoid::FindFood()` 中检测特定 `boid` 的类型，看它是天敌还是猎物，以便避开或捕食。

这些特性会增加开销吗？是的，虽然不大。当然，是否采用这样的设计取决于实现的需求。例如，在 RTS 游戏中，每队弓箭手都没有什么特殊，而在 FPS 游戏中，每个班中的成员都有其独有的特征。

#### 4. 物群的喂养

如果不同类型的 `boid` 以其他类型的 `boid` 为食，则需要控制 `boid` 的饥饿程度。除了电影 *Jaws* 中的动物外，任何动物都不是总是吃不饱的，`boid` 也不例外。

为表示这一点，老鹰和麻雀都有一个饥饿值（`hunger rating`），在每个更新周期中，该值都不断递减。当这个值为 0 时，`boid` 将感到饥饿，并开始积极地寻找猎物以填饱肚子。在我们的演示程序中，老鹰捕食麻雀，而麻雀捕食昆虫。昆虫位于食物链的最底层，不吃任何东西，但当物群中成员数量足够多时将能够繁殖。每当老鹰和麻雀捕食猎物后，将对其当前饥饿值同初始饥饿值进行比较，以判断它是否吃饱了。例如，如果老鹰开始的饥饿值为 10 点，则捕食 4 只麻雀后，它吃饱且不再捕食麻雀的可能性将为 40%。

由于作者要创建一个复杂的捕食模拟程序，遇到喜欢的食物时，老鹰和麻雀都会捕食。老鹰将试图接近它见到的麻雀，而麻雀将向最近的昆虫靠近。当然，被吃掉的 `boid` 将被删除。

#### 5. 昆虫不会绝迹

由于老鹰以麻雀为食，而麻雀靠昆虫充饥，因此昆虫位于食物链的最底层，且是生态环境中最重要的成员。如果昆虫绝迹了，其他 `boid` 也逃脱不了灭绝的命运。为防止出现这种情况，昆虫具备一种其他 `boid` 类型没有的特性——能够繁殖。为此，在 `CBoid` 类中增加了一个 `Reproduction` 参数，用于控制昆虫的繁殖。当昆虫群中的成员足够多时，将能够繁殖——每隔几秒钟就生出一个新的昆虫。

#### 6. 忙碌的晚餐

正如前面指出的，捕食是件很简单的事情。老鹰和麻雀感到饥饿时（即饥饿值降为 0 时），将寻找最近的食物，并试图接近它。老鹰捕食一只麻雀后，其饥饿值将增加 1 点，然后检查它是否吃饱了；而被吃掉的麻雀将被立刻删除。麻雀以同样的方式捕食昆虫。



由于老鹰的飞行速度通常比麻雀快，而麻雀通常比昆虫快，因此麻雀和昆虫要逃脱火顶之灾只能依靠捕食者遇到麻烦。当然，每种 boid 都将尽可能避开它见到的捕食者，这种行为会有一定的帮助。有趣的是，能够妨碍捕食者的惟一因素是我们将散布在世界中的障碍物。为避开障碍物，捕食者不得不降低速度，这让猎物有机会逃脱——总之是暂时逃脱。

### 3.11.2 有齿物群的模拟

上述工作的结果可在附带光盘中的演示程序 FlockongwithTeeth 中看到。该演示程序拥有与原有演示程序相同的控制特性，用户可能根据需要放大和缩小，开启弹性气泡 (bounding bubble)，以更好的查看 boid 的视野、避开特性和内聚距离等。

在布满大小各异的障碍物的世界中，几个老鹰（较大的三角形物体）在那里盘旋。成群的麻雀（中等大小的三角形物体）在大量的昆虫（像素大小的物体）间掠过。老鹰偶尔因饥饿而向麻雀群飞扑过去，麻雀四处逃窜以免成为老鹰的盘中餐。捕食者和猎物常常受到障碍物的妨碍，以至于无法捕获猎物或逃脱。散开的麻雀和昆虫将寻找其他劫后余生的同伴，并组成新的物群，然后整个循环再次重新开始。

演示程序最终以两种方式之一结束，结束的方式取决于配置和运气。最可能的结果是一场生态灾难——麻雀吃光了所有的昆虫或老鹰将所有的麻雀都消灭殆尽。如果昆虫灭绝了，则麻雀也将因缺乏食物而绝迹，然后老鹰也将因没有食物而绝种。如果麻雀因老鹰的捕食能力太强而灭绝，则老鹰也将跟着死光，留下一个只有昆虫和障碍物的世界。就作者的测试而言，这种情况出现的可能性较大。

另一种可能情况是，在麻雀和昆虫之间形成某种静态平衡。老鹰因运气不好捕获不到任何麻雀而灭绝。在这种情况下，麻雀将继续生存下去并以昆虫为食，但生存时间是不确定的。麻雀最终可能吃光昆虫，这样将进入前面描述的第一种情形；但也可能到达一种平衡，即昆虫的繁殖速度足够快，提供了麻雀生存下去所需的食物。

如果麻雀和老鹰也能够繁殖（演示程序不是这样，但这是一个留给读者的很简单的练习），则有可能出现任何结果。这是最逼真的配置模拟的方式，但也最难以进行正确的平衡。

### 3.11.3 局限性和可改进的地方

虽然就本文的目的而言，作者对该演示程序总体还是比较满意的，但还是有几个需要改进的地方，尤其是想将这些代码用于游戏中时。虽然在该演示程序中，只有昆虫能够繁殖，但要让老鹰和麻雀也能够繁殖很容易。另一个可改进的地方是，为形成一个更为合理的生态循环，可让昆虫也以某种东西为食，如老鹰的羽毛。

这里对视野的处理也非常基本，所有 boid 的视角都是  $360^\circ$ ，它们能够看到任何方位的东西。实际上，大多数捕食者的视觉都很敏锐，能看到很远的地方；而被猎食的动物是近视的，且看不到后面（例如狼和羊）。要使演示程序中的视觉更逼真，让老鹰能够寻找猎物而不是守株待兔也不是件困难的事，虽然限制视觉将带来大量的数学开销（视线检测和视角检测等）。

## 致谢和其他资源

这里要再次感谢 Mitre 公司的 Christopher Kline 提供的优秀的计算坡度(roll)/间距(pitch)/偏航角(yaw)的方法(这里几乎照抄了这种方法。这种方法最初是在 Christopher Kline 的 C++ *Boids* 实现中发布的,可在其网站中找到[Kline96])。读者不可不读的另一个网站是由“物群模拟技术之父”Craig Reynolds[Reynolds00]维护的,该网站中包含上百个到关于在游戏和电影中使用物群模拟技术的链接以及大量有关该主题的研究资源。作者的网站[Woodcock01]中也包含大量到关于物群模拟的文章和网页的链接。

一位读者在阅读《游戏编程精粹 1》后创建了一个小型的很棒的物群模拟程序,在该程序中可以实时地改变每种导航行为的影响力[Grub01]。这是一种很不错的方式,可以研究每种行为如何影响物群,以及它们如何一起导致突发性行为,使物群模拟技术如此有趣。

最后要介绍一本优秀的图书——Steven Levy 撰写的 *Artificial Life* [Levy93]。除讨论物群模拟技术和 boid 外,该书还概括性地介绍了人工生命(artificial life)。

### 3.11.4 参考文献

[Grub01] Grub, Tom, “Flocking Demo,” [www.riversoftavg.com/flocking.htm](http://www.riversoftavg.com/flocking.htm), March 6, 2001.

[Kline96] Kline, Christopher, maintains one of several excellent pages on flocking, together with many demos and sample code, at [www.media.mit.edu/~ckline/cornellwww/boid/boids.html](http://www.media.mit.edu/~ckline/cornellwww/boid/boids.html), August 14, 1996.

[Levy93] Levy, Steven, *Artificial Life: A Report from the Frontier Where Computers Meet Biology*, Vintage Books, 1993.

[Reynolds87] Reynolds, C. W. (1987) *Flocks, Herds, and Schools: A Distributed Behavioral Model*, in *Computer Graphics*, 21(4) (SIGGRAPH '87 Conference Proceedings) pp. 25-34.

[Reynolds00] Reynolds, Craig, maintains an extensive reference on flocking and steering behavior at [www.red3d.com/cwr/boids/](http://www.red3d.com/cwr/boids/) and has presented a wide variety of papers at various conferences discussing the progress he's made in exploring and perfecting uses of this technology, December 6, 2000.

[Woodcock00] Woodcock, Steven, “Flocking a Simple Technique for Simulating Group Behavior,” *Game Programming Gems*, Charles River Media, 2000.

[Woodcock01] Woodcock, Steven, maintains a page dedicated to game AI at [www.gameai.com](http://www.gameai.com), 2001.





## 3.12 一个用 C++ 编写的通用模糊状态机

Eric Dybsand

eric@geta.cncdsl.com

在《游戏编程精粹 1》中，由 Mason McCuskey 撰写的“视频游戏中的模糊逻辑” [McCuskey00] 以及作者撰写的“一个有限状态机类” [dybsand00] 分别介绍了模糊逻辑 (fuzzy logic) 和通用的有限状态机。本文将把这两个概念嫁接起来，编写了一个通用的模糊状态机 (Fuzzy State Machine, FuSM) C++ 类，您可以将其用于游戏中；另外本文还将简要地介绍模糊逻辑以及如何在游戏中使用模糊逻辑。

首先，简要地复习一下模糊逻辑的 FAQ 定义：

模糊逻辑是传统 (布尔) 逻辑的超集，能够处理部分真——位于完全正确和完全错误之间的真值——的概念。

因此，模糊状态机可以处于几乎 ON、差不多 OFF 或不完全 TRUE、有几分 FALSE 等状态，而不是 ON 和 OFF 或 TRUE 和 FALSE 等状态。甚至可以既 ON 且 OFF 或既 TRUE 且 FALSE，但程度不同。

对于游戏开发人员来说，这意味着什么呢？意味着非玩家角色 (Non-Player Character, NPC) 对玩家而言，不必仅仅是 MAD，而可以是差不多 MAD、部分 MAD、真正 MAD 或非常 MAD。换句话说，通过使用 FuSM (实现模糊逻辑)，游戏开发人员可以给角色赋予状态的多种程度。这还意味着游戏中的状态不必是明确和离散的 (在模糊逻辑领域，被称为明晰的 (crisp))，而可以是模糊的。其优点将在下一节讨论。

在 FuSM 内部，状态通常用 0.0~1.0 之间的实数表示，但这并非表示模糊值的惟一方式。我们可以选择任何一组数，并将其视为模糊值。继续 NPC 态度的范例 [Dybsand00]，我们来考虑如何在 FuSM 中记录 NPC 对玩家的态度中“不喜欢的成分”。可以使用 1~25 的整数来表示 NPC 对玩家的讨厌程度，26~50 之间的整数来表示 NPC 对玩家的恼怒程度，51~99 之间的整数表示 NPC 的各种愤恨程度。这样，在对玩家的每种态度中，NPC 都有不同程度的不喜欢，如讨厌、恼怒或愤恨 (见图 3.12.1)。

结束对 FuSM 的简要介绍之前，有必要澄清一个常见的关于模糊逻辑的错误认识：模糊值和概率之间没有任何关联。模糊逻辑并非一种表示概率的新方式，它表示的是集合中成员资格的程度。在概率中，值的总和必须为 1.0 才是有效的。模糊逻辑值无需满足这样的要求 (请看上述重叠的范例)。这并不意味着模糊值的总和不能是 1.0，而只是意味着要使 FuSM 有效，模糊值的总和不一定非得为 1.0。

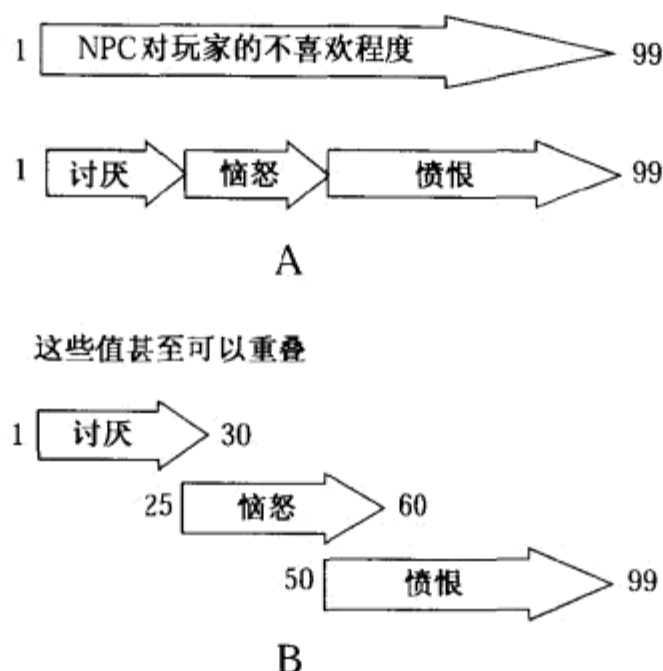


图 3.12.1 A) 表示对玩家的不喜欢态度的模糊值 B) 重叠的模糊值

### 3.12.1 为何在游戏中使用 FuSM

作者认为，在计算机游戏中使用 FuSM 的首要原因是，这是一种实现模糊逻辑的简明方式，而模糊逻辑能够深入描述用于表示游戏世界及游戏中对象之间的关系的抽象概念。

这实际上增加了玩游戏的乐趣！

您可能会问，FuSM 如何增加玩游戏的乐趣呢？FuSM 通过让 NPC 有更有趣的响应、降低 NPC 行为的可预测性，给人类玩家提供更多的选择。

这样，玩家遇到的 NPC 在受到攻击时，不会要么 MAD，要么不 MAD。相反，玩家必须应付一个有不同 MAD 程度的 NPC。这样可以为 NPC 开发更多的响应方式，让玩家见到更多的响应方式，从而增加了玩家需要考虑的因素，进而提高了玩游戏的乐趣。

将 FuSM 加入到计算机游戏中的另一个作用是提高游戏的可重玩性 (replayability)。通过增加响应方式以及玩家在特定位置遇到的情形，玩家将在类似的位置遇到不同的结果。

### 3.12.2 如何在游戏中使用 FuSM

实际上，在很多计算机游戏中已经使用了各种形式的 FuSM。

一个在计算机游戏中使用 FuSM 的例子是 NPC 或主体的生命值。主体不是要么活着要么死去 (有限状态)，而是使用一个生命值范围来反映主体从非常健康到不太健康、要死了再到死去的各种情况 (模糊状态)。另一个在计算机游戏中使用 FuSM 的例子是，在赛车游戏中控制受 AI 控制的车辆的加速或制动。在这种情况下，使用 FuSM 提供不同的加速和制动程度，而不仅仅是 THROTTLE-UP 或 THROTTLE-DOWN 和 BRAKE-ON 或 BRAKE-OFF 等有限状态。正如表示态度的例子表明的，FuSM 非常适合用于表示 NPC 的情绪状态以及它对玩家或其他 NPC 的态度。

正如前面的例子表明的，将模糊逻辑用于计算机游戏中的状态相对简单。模糊逻辑非常

适合用于结果有两个以上的决策过程，这样的过程有很多。

接下来我们将模糊逻辑加入到一个通用的 C++ FuSM 类中。

### 3.12.3 复习《游戏编程精粹 1》中的 C++ 通用有限状态机

这个用 C++ 编写的通用 FSM[Dybsand00] 由两个类组成：FSMclass 和 FSMstate。FSMclass 类对象封装了实际的有限状态机进程，记录了 FSM 的当前状态，支持各种 FSMstate 类对象的容器，提供了对状态切换过程的控制。

FSMstate 类对象封装了具体的状态，并维护输入状态数组和输出状态数组，后者指出了该对象能从当前状态切换到哪些状态。

给 FSM 的输入被提供给 FSMclass::StateTransition()，后者决定由哪个 FSMstate 对象来处理输入（根据当前状态）；然后，输入被传递给 FSMstate::GetOutput() 以获得新的输出状态。

新的输出状态被 FSMclass::StateTransition() 返回给 FSM 用户进程，并成为 FSMclass 对象的新的当前状态。因此 FSM 在响应输入时，提供的是明晰和离散的状态切换。

### 3.12.4 将通用 FSM 修改为 FuSM

要将通用 FSM 改造为 FuSM，只需修改几个地方。首先需要让 FSMclass 类对象支持多种当前状态。然后需要修改 FSMstate 类，使之支持状态程度。最后，需要修改这两个类对象中的状态切换过程，使之支持从多种状态切换到新状态的某种程度。修改过程中，我们将把 FSMclass 和 FSMstate 类重新命名为 FuSMclass 和 FuSMstate。



对于这里指出的各个类，请读者查看附带光盘中的模糊/有限状态机工程及其 ON THE CD 中的代码。

首先修改 FuSMclass，它与 FSMclass 类似，但支持多个当前状态（FuSMclass 和 FuSMstate 类的成员用粗体表示）。这种功能是由 m\_list 成员 FuzzyState 提供的，该成员是一个 STL 链表对象，其中包含指向 FuSMstate 对象的指针。对于任何活动的模糊状态对象（取决于当前的输入值），都将指向它的指针存储到该链表中。这样，便能够支持多个当前状态。和 FSMclass 一样，FuSMclass 也维护一个 STL 映射表对象（m\_map 成员 fuzzyState\_Map），该对象中包含指向 FuSMclass 可能考虑的所有 FuSMstate 对象的指针。

接下来开发一个名为 GetNextFuzzyStateMember() 的访问成员函数，它提供存取 FuSMclass 对象的服务。成员函数 GetNextFuzzyStateMember() 维护一个指针，该指针指向 m\_list 成员 FuzzyState\_List 中的下一个 FuSMstate 指针，以便 FuSMclass 外面的进程能够访问所有的当前活动状态。因此，这种服务是程序如何存取当前活动状态。通过不断调用 GetNextFuzzyStateMember()，直到它返回一个 NULL 指针，程序能够获悉所有的活动模糊状态。

接下来修改 FuSMstate 类，以支持各种程度的成员资格。这种支持是通过新增成员变量 int m\_iLowRange 和 int m\_iHighRange 来提供的。出于简化的目的，该设计将模糊成员资格范围视为正整数，可以很容易地对其进行修改，将模糊成员资格范围视为一组实数。出于方便

的考虑，还维护了 FuSMstate 对象的另外两个属性：FuSMstate 对象的成员资格值（int m\_iValueOfMemberShip）和成员资格程度（int m\_iDegreeOfMemberShip）。

注意，无限状态对象（FSMstate）和新的模糊状态对象（FuSMstate）之间的最主要的区别是，不再需要状态切换数组。这是因为在模糊逻辑中，同时可以有一个或多个状态；而在有限逻辑中，同时只能有一种状态。

最后，修改 FuSMclass 和 FuSMstate 对象中的状态切换过程，以支持多个当前状态以及给定状态中各种程度的成员资格。对于 FuSMclass，这意味着需要修改 StateTransition() 以处理包含所有可能状态的 m\_map 成员 FuzzyState\_Map，让每个 FuSMstate 对象都能够根据累积的输入值（FuSMclass 中的 int 成员 m\_iCurrentInput）进行状态切换。指向进行状态切换的 FuSMstate 对象的指针被存储在 m\_list 成员 FuzzyState\_List 中，从而表明这些 FuSMstate 对象是当前活动状态。

对于 FuSMstate 类对象，需要用一个新的状态切换函数替换（FSM 中的）FSMstate::GetOutput() 函数。新的成员函数 FuSMstate::DoTransition() 接受 FuSMclass 维护的输入值，并考虑该输入值表示的成员资格程度。如果模糊状态中的成员资格存在，则该函数返回 TRUE，并维护该成员资格的状态，供以后访问。



至此，修改工作便完成了。有关更详细的细节和代码清单，请参阅附带光盘中的 FuSM 工程。

### 3.12.5 在游戏中使用模糊逻辑

---

参考 FuSMclass 和 FuSMstate 类，可以提高游戏的模糊性。这样做可以让玩家在玩游戏时有更丰富的体验，同时您将更深入地理解如何应用游戏开发人员可用的、更灵活的人工智能工具。

### 3.12.6 参考文献

---

[Dybsand00] Dybsand, Eric, "A Generic Finite State Machine in C++," *Game Programming Gems*, Charles River Media, 2000.

[FAQ97] "What is fuzzy logic?" FAQ: Fuzzy Logic and Fuzzy Expert Systems 1/1 Monthly Posting, [www.faqs.org/faqs/fuzzy-logic/part1/](http://www.faqs.org/faqs/fuzzy-logic/part1/), 1997.

[McCuskey00] McCuskey, Mason, "Fuzzy Logic for Video Games," *Game Programming Gems*, Charles River Media, 2000.



### 3.13 避免模糊系统中的组合激增

Michael Zarozinski, Louder Than A Bomb! Software  
michaelz@louderthanabomb.com

归根结底，模糊逻辑不过是一系列的 if-then 语句。使用模糊逻辑时，最大的问题之一是，随着模糊状态数的增加，if-then 语句的数目将呈几何级数增长。这被称为组合激增（combinatorial explosion），将导致模糊系统缓慢、混乱且难以维护。在游戏中，速度至关重要，而组合激增将导致无法使用模糊逻辑。

有关模糊逻辑的简介，请参阅《游戏编程精粹 1》中 Mason McCuskey 撰写的《视频游戏中的模糊状态》一文[McCuskey00]。本文将介绍一些定义，因为模糊逻辑方面的术语很混乱。

- **变量 (variable)**: 模糊变量指的是诸如温度、距离或生命值等概念；
- **集合 (set)**: 在传统的逻辑中，集合是明晰的 (crisp)；要么 100% 属于某个集合，要么不属于。高个子集合由所有身高超过 6 英尺的人组成，任何身高低于 6 英尺的人都属于“矮个子”（更准确地说，是不属于“高个子”）。模糊逻辑允许集合是模糊的，因此身高超过 6 英尺的人可能有 100% 的“高个子”资格，也可能有 20% 的“中等个子”资格。

#### 3.13.1 问题

表 3.13.1 说明了随着系统中变量和/或集合的增加，组合激增带来的影响。如果每次都需要检查每种可能的规则，规则数呈几何级数增长可导致任何系统崩溃。

表 3.13.1 组合激增的影响		
变量数	每个变量的集合数	规则数
2	5	$5^2 = 25$
3	5	$5^3 = 125$
4	5	$5^4 = 625$
5	5	$5^5 = 3125$
6	5	$5^6 = 15625$
7	5	$5^7 = 78125$
8	5	$5^8 = 390625$
9	5	$5^9 = 1953125$
10	5	$5^{10} = 9765625$



3.13.2 解决方案

波音公司的工程师 Willian E. Combs 发明了一种将上述几何级数增长转换为线性增长的方法，这种方法被称为“Combs 方法”。这样，对于包含 10 变量，每个变量有 5 个集合的系统，将只有 50 个规则，而不是 9 765 625 个。

注意，Combs 方法并非一种将已有的 if-then 规则转换为线性系统的算法。从一开始，您就应创建适应 Combs 方法的规则。

如果读者对 Combs 方法背后的理论感兴趣，请参阅本文最后的证明。

1. 真实世界

为将该理论用于真实世界，我们来看一个计算游戏中个体的进攻性的系统。这里只考虑一场有 3 个变量的一对一的战斗，而忽略周围的个体（友军和敌人）。这 3 个变量是：

- 自身的生命值（health）；
- 敌人的生命值；
- 敌我之间的距离。

变量“生命值”有 3 个集合：要死了、良和优；

“距离”变量也有 3 个集合：近、中、远；

最后，输出（攻击性）也有 3 个集合：逃跑、且战且退、全力进攻。

2. 传统的模糊逻辑规则

如果使用传统的模糊逻辑系统，则首先需要以电子表格的方式创建规则，如表 3.13.2 所示。

表 3.13.2 一些传统的模糊逻辑规则

自身的生命值	敌人的生命值	距 离	攻 击 性
优	优	近	且战且退
优	优	中	且战且退
优	优	远	全力进攻
优	良	近	且战且退
优	良	中	全力进攻
优	良	远	全力进攻
优	要死了	近	全力进攻
优	要死了	中	全力进攻
优	要死了	远	全力进攻
		.	
		.	
		.	
良	良	近	且战且退
良	要死了	近	且战且退



续表

自身的生命值	敌人的生命值	距 离	攻 击 性
		.	
		.	
		.	
要死了	优	近	逃跑
要死了	优	中	逃跑
要死了	优	远	且战且退

注意，表 3.13.2 只列出了 27 个可能规则中的 14 个。虽然这样的范例很容易管理，但组合激增将很快开始起作用。在游戏中，可能还需要考虑其他变量，如可能加入战斗的友军和敌军的相对生命值。增加这两个变量后（变量总数为 5 个），上表中包含的规则将从 27 个增加到 243 个。这很快将失控。幸运的是，为处理 5 个变量，Combs 方法只需要 15 条规则。

3. 关于模糊逻辑规则的 Combs 方法

在传统系统中创建规则时，我们考虑的是输入集合的组合同输出之间的关系。使用 Combs 方法时，考虑的是各个集合同输出之间的关系，每次为一个变量创建规则（见表 3.13.3）。

在使用 Combs 方法的系统中，建议所有输入变量的集合数同输出变量的集合数相同。并非一定要这样做，但它让每个输出集合与每个输入变量的集合呈一一对应的关系。

表 3.13.3 每个输入变量的集合同输出的关系

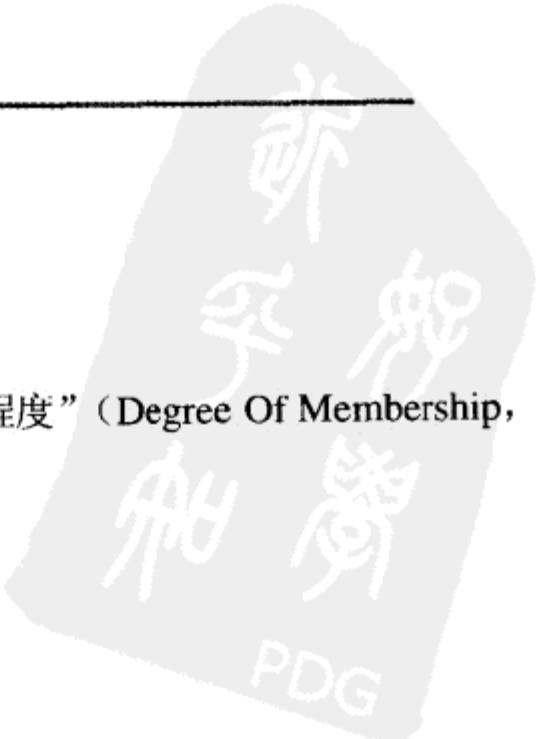
自身的生命值	进 攻 性
优	全力进攻
良	且战且退
要死了	逃跑
敌人的生命值	进攻性
优	逃跑
良	且战且退
要死了	全力进攻
距离	进攻性
近	且战且退
中	且战且退
远	全力进攻

3.13.3 范例

为测试该系统，我们将使用下面的值：

- 自身的生命值：76.88；
- 敌人的生命值：20.1；
- 距离：8.54。

图 3.13.1~3.13.3 是系统输入值的“成员资格程度”（Degree Of Membership, DOM）。注意，一个变量的 DOM 的总和可以不为 100%。



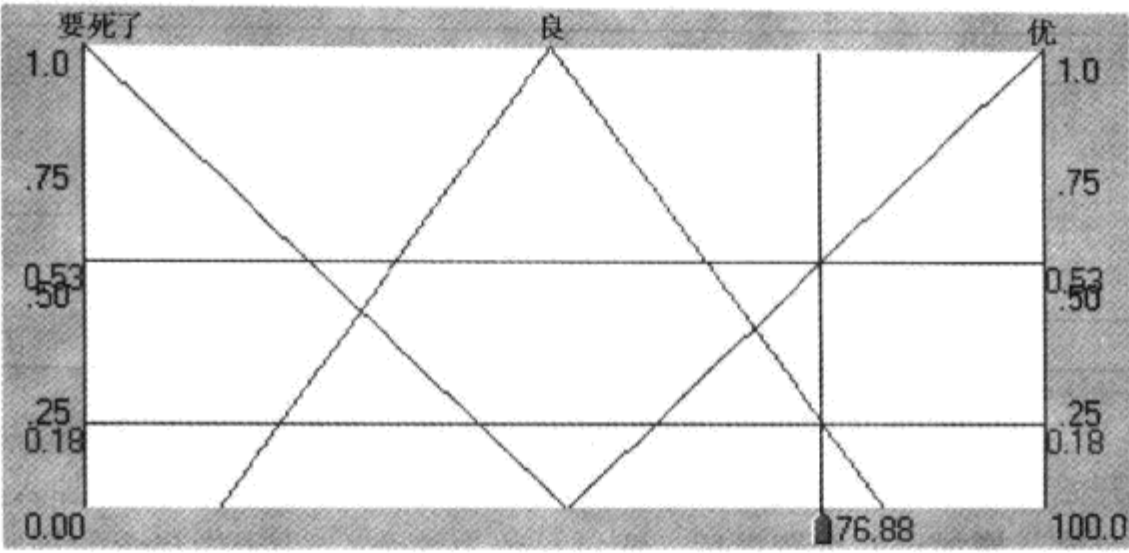


图 3.13.1 自身生命值 76.88 的 DOM: 要死了 0%、良 18%、优 53%

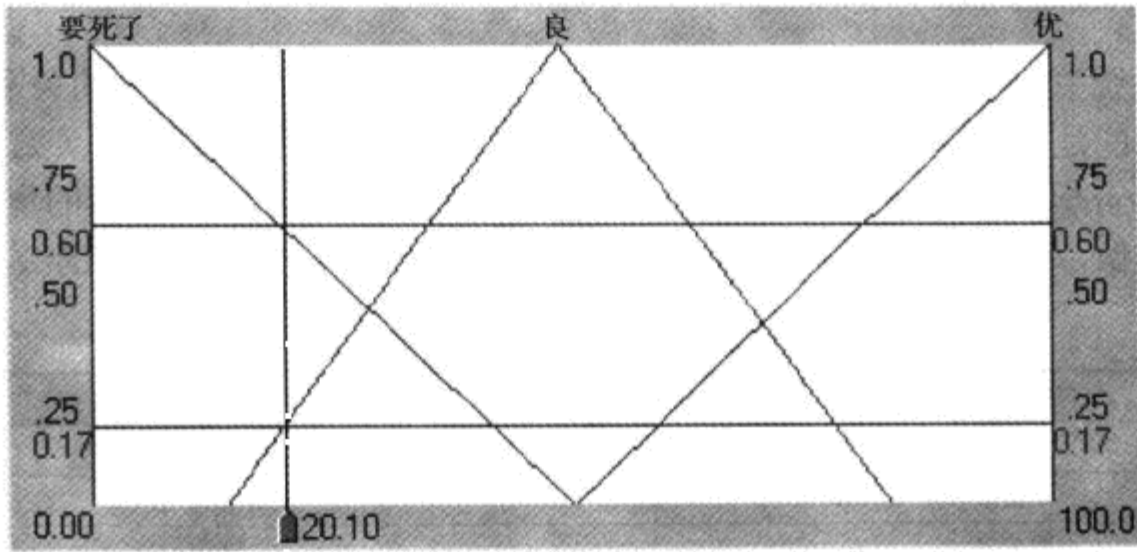


图 3.13.2 敌人生命值 20.1 的 DOM: 要死了 60%、良 17%、优 0%

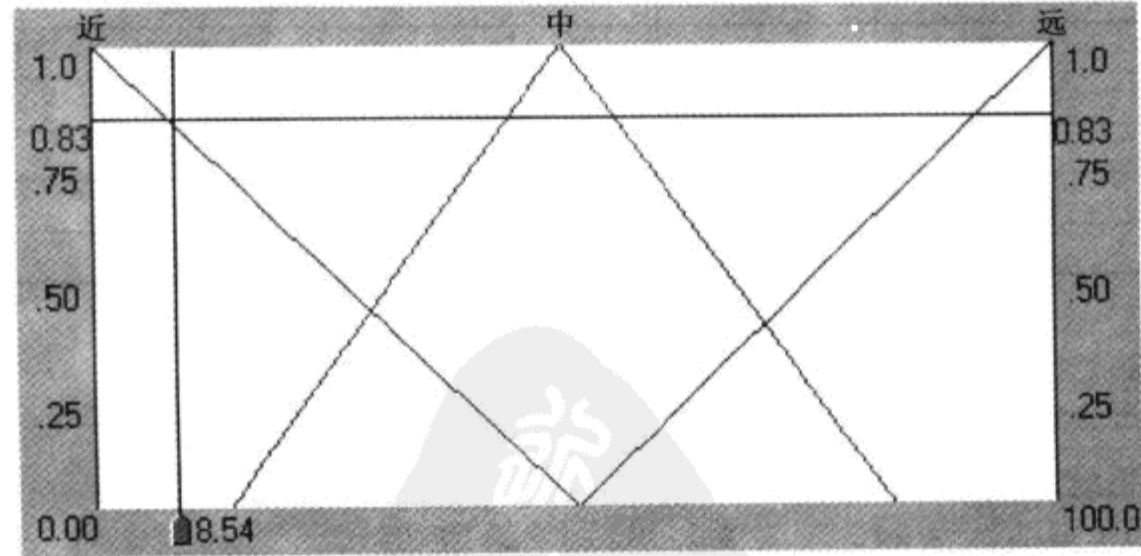


图 3.13.3 距离值 8.54 的 DOM: 近 83%、中 0%、远 0%

图 3.13.1~3.13.4 来自模糊逻辑编辑器“Spark!”，该编辑器让您能够以可视化的方式创建模糊逻辑系统、将其集成到游戏中并实时地修改 AI 而无需重新编译。

1. 传统系统范例

使用前面创建的规则时，在传统系统中将触发表 3.13.4 中列出的规则。

将对输入集合的 DOM 执行 AND 运算，以获得输出集合的 DOM。AND 运算相当于取三个输入值中最小的一个。

我们使用的逆模糊化（defuzzification）方法——质心（center of mass）——根据输出集合的最大值找出输出集合的质心（见图 3.13.4）。

表 3.13.4 在传统系统中被触发的规则

自身的生命值	敌人的生命值	距 离	攻 击 性
优 (53%)	良(17%)	近(83%)	且战且退(17%)
优 (53%)	要死了 (60%)	近(83%)	全力进攻(53%)
良(18%)	良(17%)	近(83%)	且战且退 (17%)
良(18%)	要死了 (60%)	近(83%)	且战且退(18%)

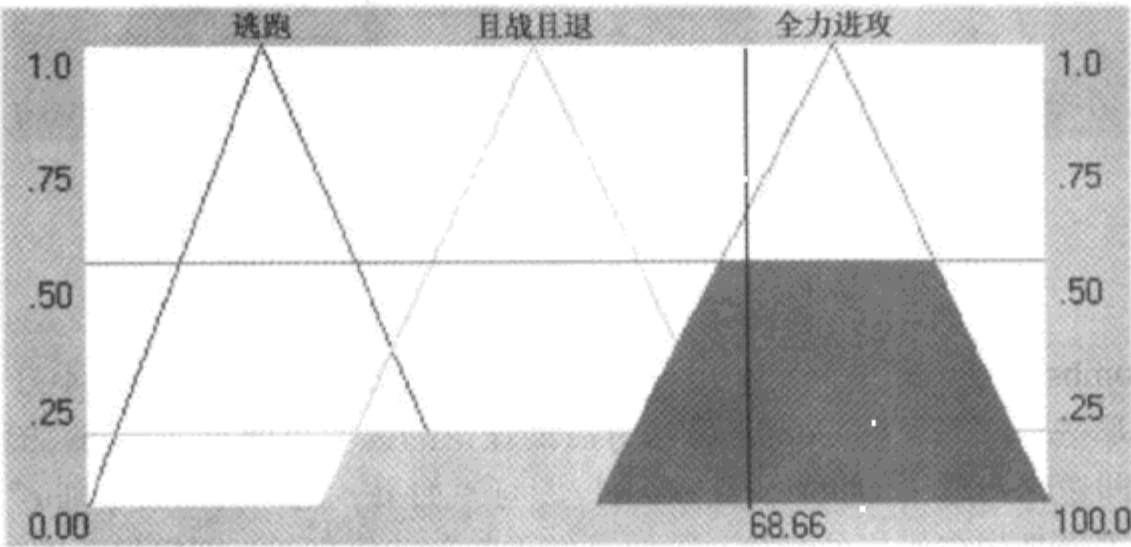


图 3.13.4 传统系统的输出：且战且退 18%、全力进攻 53%、进攻性 68.66

2. Combs 方法的范例

使用同样的输入时，Combs 方法的规则如表 3.13.5 所示。

表 3.13.5 Combs 方法系统的输出：且战且退 83%、全力进攻 53%、攻击性 60.39

自身的生命值	攻 击 性
优(53%)	全力进攻 (53%)
良(18%)	且战且退 (18%)
敌人的生命值	攻击性
良(17%)	且战且退 (17%)
要死了(60%)	且战且退 (60%)
距离	攻击性
近(83%)	且战且退 (83%)

Combs 方法的结果为 60.39，与传统方法的结果（68.66）不同，但我们本来就没有期望它们相同，因为使用的推理方法不同。Combs 方法对值执行 OR 运算，即取得其中最大的值（见图 3.13.5）。传统的模糊逻辑对值执行 AND 运算，即取得其中的最小值；因此且战且退的

输出集合不同。

注意，在这个例子中，如果取输出集合的最小值（没有规定不能这样做），则结果将与传统的模糊逻辑方法相同。结果取决于在 Combs 方法中选择的规则。由于并不存在将传统模糊逻辑规则转换为 Combs 方法的算法，因此不能说仅需取最小值便总是能够得到与传统方法相同的结果。

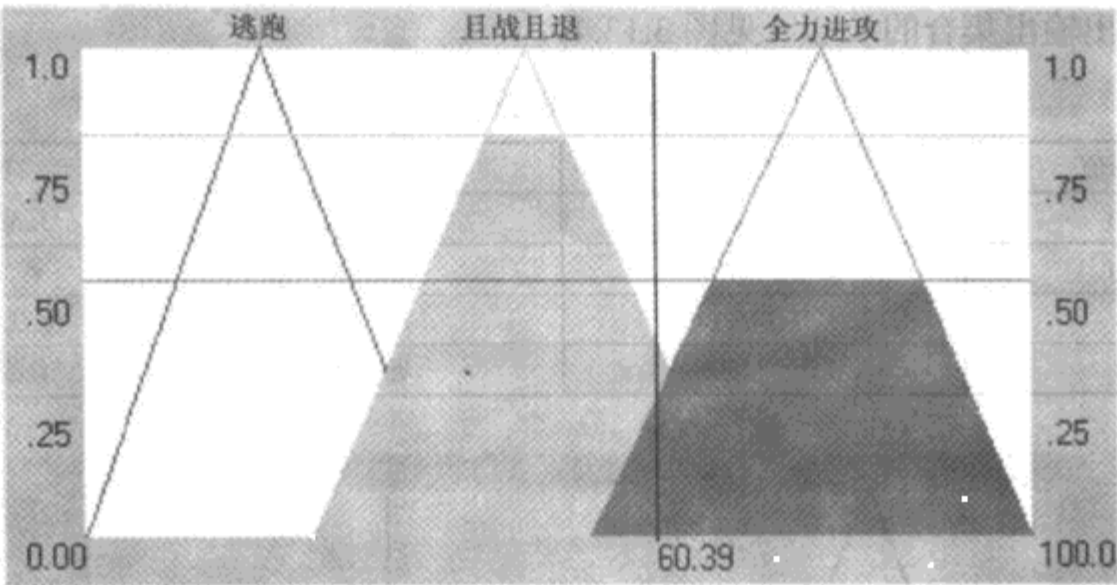


图 3.13.5 Combs 方法的输出：且战且退 83%、全力进攻 53%、攻击性 60.39

3. 证明

要使用 Combs 方法，并不一定非得理解该方法为何管用。形式逻辑并不好理解，因此这里的证明是为那些想深入理解 Combs 方法背后的理论的读者提供的。Combs 方法是基于这样一个事实的，即逻辑命题“如果 (p 且 q) 则 r”与“如果 p 则 r 或如果 q 则 r”等价。

由于模糊逻辑是形式逻辑的超集，因此我们不能忽略模糊性，则仅仅证明 Combs 方法与传统方法等价。在表 3.13.6 中，p 和 q 是前提，r 是推论。前提是诸如“如果 Jim 的个子很高”或“如果 Jim 身体健康”等陈述；而推论是可能的结果，如“Jim 能够打篮球”。

除了 x 则 y 类型的子句，这里的证明很简单。这些子句是标准的形式逻辑命题，但它们的真值表不太好理解，尤其是当 x 和 y 为假但命题为真。有关可帮助阐述该命题的真值表的例子，请参阅[aiGuru01]。

表 3.13.6 前提和推论

p	q	p 且 q	r	如果 p 且 q 则 r	如果 p 则 r	如果 q 则 r	如果 p 则 r 或如果 q 则 r
T	T	T	T	T	T	T	T
T	T	T	F	F	F	F	F
T	F	F	T	T	T	T	T
T	F	F	F	T	F	T	T
F	T	F	T	T	T	T	T
F	T	F	F	T	T	F	T
F	F	F	T	T	T	T	T
F	F	F	F	T	T	T	T

如果要以可视化的方式进行证明，请看图 3.13.6 和 3.13.7 中的维恩图。由于维恩图只能

表示 AND、OR 和 NOT 关系，因此实质蕴含做了下述转换：

- 传统逻辑：“如果  $p$  且  $q$  则  $r$ ”与“ $(\text{非 } p \text{ 且 } q) \text{ 或 } r$ ”等价；
- Combs 方法：“如果  $p$  则  $q$  或如果  $q$  则  $r$ ”与“ $(\text{非 } p \text{ 或 } r) \text{ 或 } (\text{非 } q \text{ 或 } r)$ ”等价。

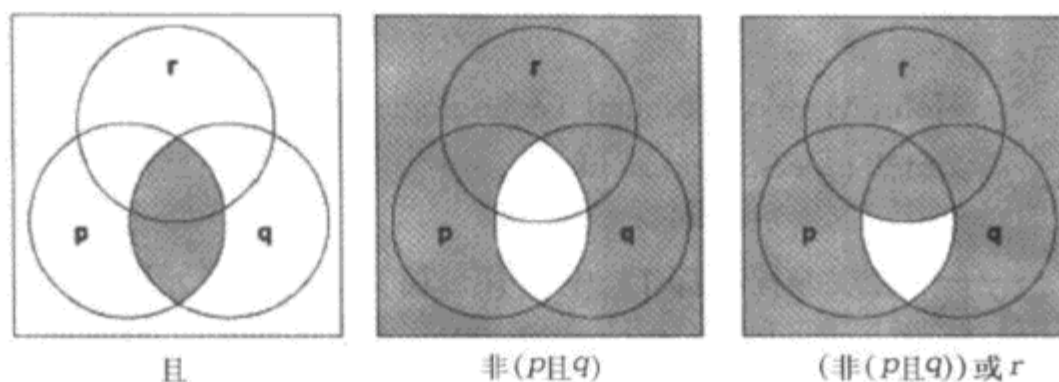


图 3.13.6 传统逻辑的维恩图

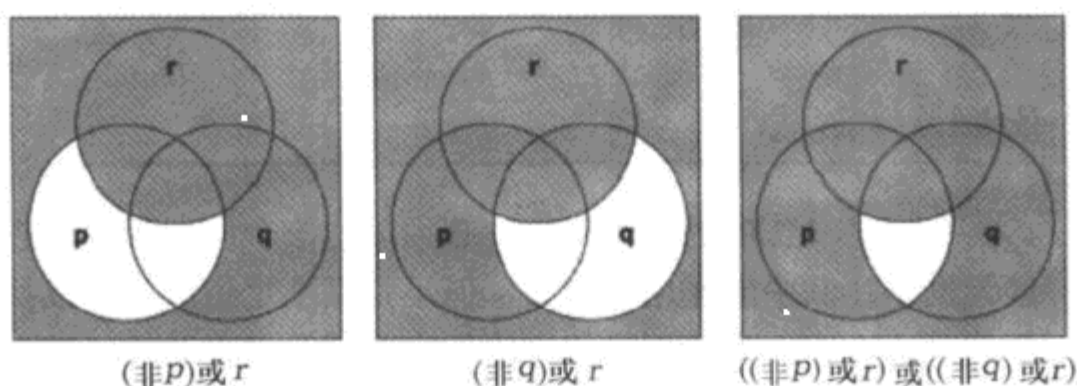


图 3.13.7 Combs 方法的维恩图

### 3.13.4 结论

如果必须要在速度和智能游戏 AI 之间选择的话，被选择的几乎总是速度（基于回的游戏除外）。Combs 方法让我们能够创建这样的 AI，即不断速度快，而且足够复杂，使得行为丰富、逼真，进而给玩家提供引人入胜的体验。

以后遇到需要创建复杂的行为系统的时候，请尝试使用 Combs 方法。您可能会惊奇地发现，使用如此少的规则就可能得到如此丰富的行为。

### 3.13.5 参考文献

[aiGuru01] Zarozinski, Michael, “if  $p$  then  $q$ ?” available online at [www.aiGuru.com/logic/if\\_p\\_then\\_q.htm](http://www.aiGuru.com/logic/if_p_then_q.htm), March 5, 2001.

[Andrews97] Andrews, James E., “Taming Complexity in Large-Scale Fuzzy Systems,” PC AI (May/June 1997): pp.39–42.

[Combs99] Combs, William E., *The Fuzzy Systems Handbook 2nd Ed*, Academic Press, 1999.

[McCuskey00] McCuskey, Mason, “Fuzzy Logic for Video Games,” *Game Programming Gems*, Charles River Media, 2000.

## 3.14 一个在游戏中使用神经网络的例子

---

John Manslow, Neural Technologies Limited  
jfm96r@ecs.soton.ac.uk

在《游戏编程精粹 1》中，有一篇文章对神经网络（neural network）做了全面的概述[Lamothe00]。本文作为该文的补充，将介绍一个具体的应用范例——多玩家感知器（MultiLayer Perceptron, MLP），这是神经网络在游戏中最大的用途之一，也是使用最广泛的。为此，将描述如何识别可通过 MLP 解决的问题，并阐述实现解决方案的步骤。

### 3.14.1 游戏

---

为提供 MLP 的应用范例，必须创建一个应用程序，它足够简单，神经网络在其中的作用很清晰，但不是微不足道的，可以找到显而易见的替代解决方案。为此创建的应用程序是一个坦克游戏，其中两辆坦克被放在随机生成的地形的侧视图中，左边的坦克由玩家控制，而右边的坦克由计算机控制。

坦克轮流地向对方射击，首先命中对方者取胜。每辆坦克通过调整其炮筒的倾斜角来瞄准——这非常困难，因为坦克前进时，其发射炮弹的速度将降低（由于阻力），并受到风力的影响（出于简化的目的，假设在发射炮弹期间，风速和方向保持恒定）。因此，AI 面临的主要挑战是，如何设置坦克炮筒的倾斜角，以命中敌方坦克。本文将介绍如何对 MLP 进行训练，以解决这种问题。

### 3.14.2 多玩家感知器

---

MLP 是一种神经网络，在 20 世纪 80 年代的中期，随着人们发现一种非常有效的训练方式，MLP 开始流行。从那时起，这种技术日益流行，现已成为业界应用最为广泛的神经网络体系结构之一，并在大量的游戏中被采用（如 Codemasters 的 *Colin McRae Rally 2.0* 中）。

虽然现在在学术界，更先进的技术并不鲜见，但 MLP 仍在流行，这是因为它是最容易理解、编码最容易、应用最容易的技术之一，而且它即使被经验不那么丰富的用户使用也能提供健壮的性能。因此对于不熟悉神经网络的人来说，MLP 是一个不错的着手点；而对经验比较丰富的用户来说，它又是一个功能强大的工具。



诸如 MLP 等神经网络 (Neural network) 实际上只不过是一个复杂的非线性函数, 该函数接受大量可调整的参数, 可以通过修改这些参数来控制其形态 (shape)。训练神经网络只不过是调整其参数, 使其表示的函数有期望的形态。虽然也可以采用多项式和样条线来解决与 MLP 类似的问题, 但 MLP 的结构使其非常健壮。

函数训练成的形态是由输入-输出样本对指出的, 因此神经网络训练只不过是曲线拟合而已——调整网络中的参数, 使其大概与样本吻合。学习过自然科学课程的读者应熟悉这一过程, 在自然科学中, 常常需要根据通过实验数据绘制出一条平滑的曲线。

由于神经网络表示的通常是非常复杂的方程, 因此常常用有向图 (directed graph) 来表示神经网络。例如, 方程 3.14.1 是一个表示 MLP 的公式, 该 MLP 的输出  $y$  是线性的, 同时接受  $N$  个输入—— $x_1$  到  $x_N$ , 并包含  $M$  个隐藏的神经元。图 3.14.1 是该 MLP 的图形表示。虽然 MLP 是用有向图表示的, 但通常省略网络图中的箭头, 因为信息总是沿输入到输出的方向流动。

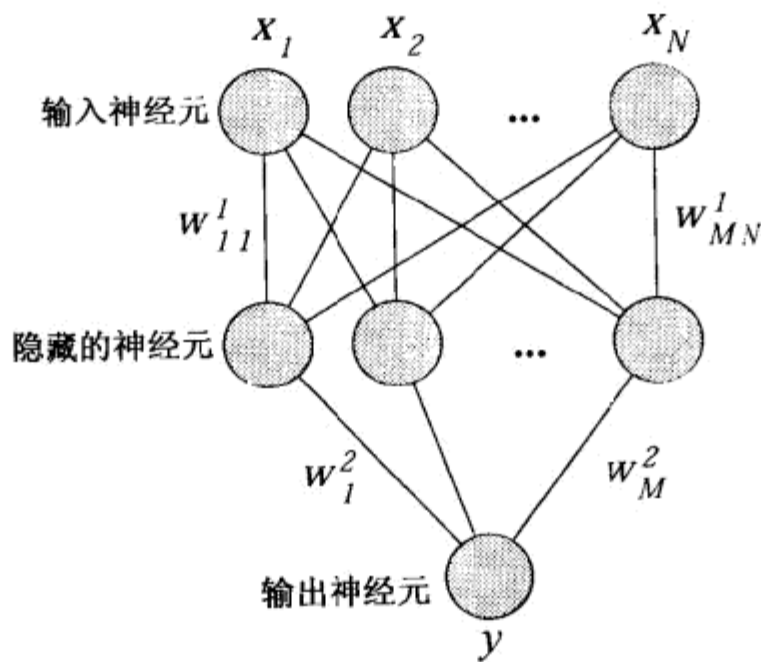


图 3.14.1 方程 3.14.1 对应的 MLP 的图形表示, 其中的圆表示可以变化的输入神经元和隐藏神经元的数目

$$y = b^2 + \sum_{m=1}^M \left[ w_m^2 \frac{1}{1 + \exp \left[ -b_m^1 - \sum_{n=1}^N w_{mn}^1 x_n \right]} \right]$$

网络中可调整的参数是  $w$  (权重) 和  $b$  (偏压)。权重将神经元连接起来, 而偏压刺激或抑制神经元, 即使网络中没有其他活动。图 3.14.1 中没有画出偏压, 因为偏压可被视为神经元的内部组分。有关神经网络的结构详细介绍, 请参阅[LaMothe00]和[Haykin94], 这里不再介绍。

就本文而言, 只需知道 MLP 表示方程 3.14.1 所示的函数, 该函数将被用来计算 AI 坦克的炮筒的倾斜度以便能够命中玩家的坦克即可。方程 3.14.1 中的  $w$  和  $b$  是可调整的参数, 我们可以使用这些参数来使 MLP 与一组样本吻合, 这组样本指出了应如何设置炮筒的倾斜角。

这些参数的值是通过曲线拟合程序来确定的，这一过程也被称为训练。

因为，为训练 MLP，需要一组由输入—输出对组成的样本，这些样本是一些这样的例子，即为命中玩家的坦克，应如何设置 AI 坦克的炮筒。显然，根据我们要解决的问题，我们希望神经网络的输出为 AI 坦克的炮筒的倾斜角，但输入应是什么呢？这常常是一个复杂的问题，将在下一节讨论。

### 3.14.3 选择输入

收集训练网络所需的输入—输出样本之前，必须确定 MLP 需要什么样的输入。显然，输入必须包含 MLP 计算正确的输出所需的信息，在这个例子中，是命中玩家的坦克所需的炮筒倾斜角。

在实践中，输入的选择常常是困难的，因为可从游戏世界中提取的信息通常很多，而要解决的问题常常太复杂或对其认识不深刻，无法确定什么样的信息是有用的。在这种情况下，可能需要花大量的时间，使用不同的输入变量组合，对很多不同的网络进行训练，以确定哪种组合的效果最佳。为使找到良好组合所需的工作量尽可能小，应遵循下述指导原则：

- 使用关于要解决的问题的先验知识，推测游戏世界中的哪些信息可能与问题相关。如果认为网络输出应取决于变量的特定函数，则将该函数加入到网络输入集合中。
- 由简单变量组成的函数导出的抽象变量提供的信息通常比前者更精确。例如，在策略游戏中，单个敌军兵力指示器可能比多个与敌军兵力的不同方面相关的指示器更有用。
- 使用尽可能提供关于游戏世界不同信息的变量，这样通过更少输入便能提供同样多的信息。例如，同时以极坐标和迪卡尔坐标的方式提供玩家坦克和 AI 坦克的相对位置不会有任何好处，因为它们包含的信息相同。
- 尝试使用组合搜索算法来得到良好的输入组合。例如，预先选择大量的候选输入，不断地在网络中增加单个输入，每个阶段增加最能提高性能的那个输入。这种技术需要的人工干预很少，但速度很慢，可能无法找到最佳组合。

虽然允许网络使用所有可能相关的变量以避免费力的输入选择工作很有诱惑力，但这样的网络的行为很可能可能不可预测，并且很糟糕。输入选择是开发神经网络应用程序中劳动密集度最高的部分，找到一小组包含大量相关信息的输入是成功的关键。

幸运的是，对于计算 AI 坦克炮筒的倾斜角，我们的先验知识足以准确地断定需要哪些信息——坦克之间的距离（这可以用两个输入来表示： $x$  距离和  $y$  距离）以及风速和方向（可用一个有符号输入来表示）。炮弹受到的阻力是固定的，因此无需将阻力的大小告知网络，但网络将在训练期间了解阻力的影响。确定需要使用的输入后，便可以收集用于训练的样本。

### 3.14.4 收集数据

为生成表明如何设置 AI 坦克炮筒的倾斜角样本，最简单也是最显而易见的方式是让玩家来控制 AI 坦克，并记录选择的输入变量（坦克的相对位置和风速）以及每次命中时 AI 坦克的炮筒的倾斜角。这一过程将不断重复下去，直到有足够多的样本供训练使用。由于人类玩家通常发射 5 次才能命令敌方坦克一次，而训练需要数百甚至数千个样本，因此收集数据

的时间可能很长。

为生成本文使用的样本，通过执行简单的随机搜索来找出正确的倾斜角，从而将数据收集过程自动化。这是通过将 AI 坦克的炮筒的倾斜角设置为一个随机的值，直到命中目标来完成的。命中目标后，便记录此时坦克的相对位置、风速和炮筒的倾斜角。然后重复这一过程——随机地放置坦克和设置风速。

虽然这种方法的效率非常低，但它无需人工干预，将不断运行直到生成了足够大的数据集。另外，这种游戏世界模拟无需以正常速度进行，同时可以通过关闭渲染（如果这样做不会改变行为的话）来提高速度。只需一晚，这就能够实现 1200 次命中目标，从而生成大约 1200 个样本。

收集数据时，需要询问的一个重要问题是“应收集多少数据？”。不幸的是，这个问题没有简单的答案，因为这取决于要解决的问题的复杂程度。在实践中，如果隐藏的神经元很少（本文使用 10 个，这通常是一个很大的数目），则对于使用  $I$  个输入的网络， $10^I$  个训练样本就能取得不错的效果。

开发神经元网络解决方案的一个重大风险是，最后调整游戏引擎时，常常会对玩游戏的方式进行一些微调，如果这改变了游戏世界的行为，将导致神经网络的表现非常糟糕。为解决这种问题，可重复数据收集和训练过程，为网络生成一组新的参数。如果游戏行为的变化不大，则通常无需重新选择输入。确定网络需要控制什么及其所需的输入，并收集了一些告诉网络如何做的训练数据后，便可以开始训练网络。

### 3.14.5 训练 MLP

正如前面指出的，MLP 是一个非线性函数，它通过调整其参数来与一系列的样本相吻合。训练的方法如下：使用一种优化算法来确定参数的值，这些参数值使得对于给定的输入，MLP 生成的输出与实际输出之间的误差最小。最常见的误差度量是均方误差，它是 MLP 的输出和样本中相应输出之间的差的平方和再除以样本数。

虽然优化算法“梯度下降”（gradient descent）常被用于使 MLP 与训练样本相吻合，但这里将采用一种很少被使用的技术——扰动搜索（perturbation search），因为它更容易编码、更容易理解，应用起来也更容易（例如，它确保是稳定的）。另外，扰动搜索不需要梯度信息，这样：

- 便于对各种网络结构、非线性特性和误差函数进行试验；
- 消除了由于错误地计算梯度信息或在网络结构中传递这种信息引发的 bug；
- 允许直接对网络的整数版本（针对低端平台）进行优化，避免了将浮点网络转换为整数网络引发的各种不合逻辑的行为；
- 允许网络中包含不连续的函数。

基本的扰动搜索可概要如下：测量 MLP 的性能。通过给每个 MLP 的参数加上一个随机噪声来扰动它们，然后重新测量 MLP 的性能。如果性能恶化，则将参数恢复到原来的值。不断重复这一过程，直到满足某种结束条件。

由于所有可用于训练 MLP 的技术都将无限地改进 MLP 的性能（如果能不断改进该多好），因此将就何时结束训练做出某种决策。为此，必须在游戏中定期地对 MLP 的性能进行

评估，并在性能足够好或进一步的训练无法改进性能时结束训练。

在游戏中评估 MLP 的性能时，必须以玩家实际玩该游戏的典型方式来运行游戏。这可确保测试时 MLP 面临的环境与游戏发行后 MLP 将遇到的环境类似，从而确保测试结果能够说明 MLP 在实际产品中的效果。如果游戏中 MLP 的性能无法达到有用的程度，则应考虑下述原因：

- 优化算法找到了局部最小值[Bishop95]。尝试使用一组随机的参数重新开始训练；
- 输入样本没有包含足够的关于相应输出的信息，网络无法再现输出。重新选择输入，找到包含更多相关信息的新输入。
- 网络太简单，无法认识到样本数据中输入和输出之间的关系。考虑对输入进行变换以简化输入与输出之间的关系，或者增加网络中隐藏的神经元数目（但应确保这样的神经元非常少）。
- 样本不能代表网络在游戏中遇到的环境。收集训练样本后，不要更改游戏世界的行为；样本中必须包含网络在游戏中将遇到的各项因素，且这些因素所占的比例是合适的。

#### 计算方面的问题

由于游戏 AI 运行在 CPU 时间短缺的环境中，因此必须考虑神经元网络的计算成本。不幸的是，训练 MLP 是处理器密集型的，这使得 MLP 非常不适合在游戏中学习，在大多数情况下，几乎总是可以采用更简单得多的学习机制。训练过 MLP 后，计算其输出需要的处理器时间非常少，尤其是当所有的内部量都可以用整数进行模拟且用查找表替代非线性函数时。这种优化使得 MLP 可以多种实时方式被用于 PC 和游戏站中。

### 3.14.6 结果

按前面各节概括的步骤进行，我们创建了一个使用 3 个输入的 MLP，其中两个是玩家坦克和 AI 坦克之间的相对位置的迪卡尔表示，另一个是风速。我们收集了 1207 个成功命中目标的样本，并在配置了 500MHz Intel 赛扬处理器的 PC 上对 MLP 训练了大约两个半小时。至此，训练便结束了，因为 MLP 的性能满足了要求，在游戏中的命中率为 98%。

### 3.14.7 结论



本文描述了创建用于简单坦克游戏中的神经元网络 AI 的步骤，附带光盘中提供了该坦克游戏。有兴趣的读者可通过参考文献[Haykin94]和[Bishop95]深入了解诸如输入选择和过适应（overfitting）等主题，这些文献对此做了更详细、更全面的介绍。最后，任何东西都无法代替实践经验——请对附带光盘中的 MLP 进行试验，并将其用于解决您自己的问题。按这里概述的开发过程进行，您将发现神经元网络是一个灵活、功能强大的工具。

### 3.14.8 参考文献

---

[Bishop95] Bishop C. M., *Neural Networks for Pattern Recognition*, Oxford University Press Inc., 1995.

[Haykin94] Haykin S., *Neural Networks: A Comprehensive Foundation*, Macmillan College Publishing Company, 1994.

[LaMothe00], LaMothe, A., *Game Programming Gems*, Edited by DeLoura, M., Charles River Media, Inc, 2000.

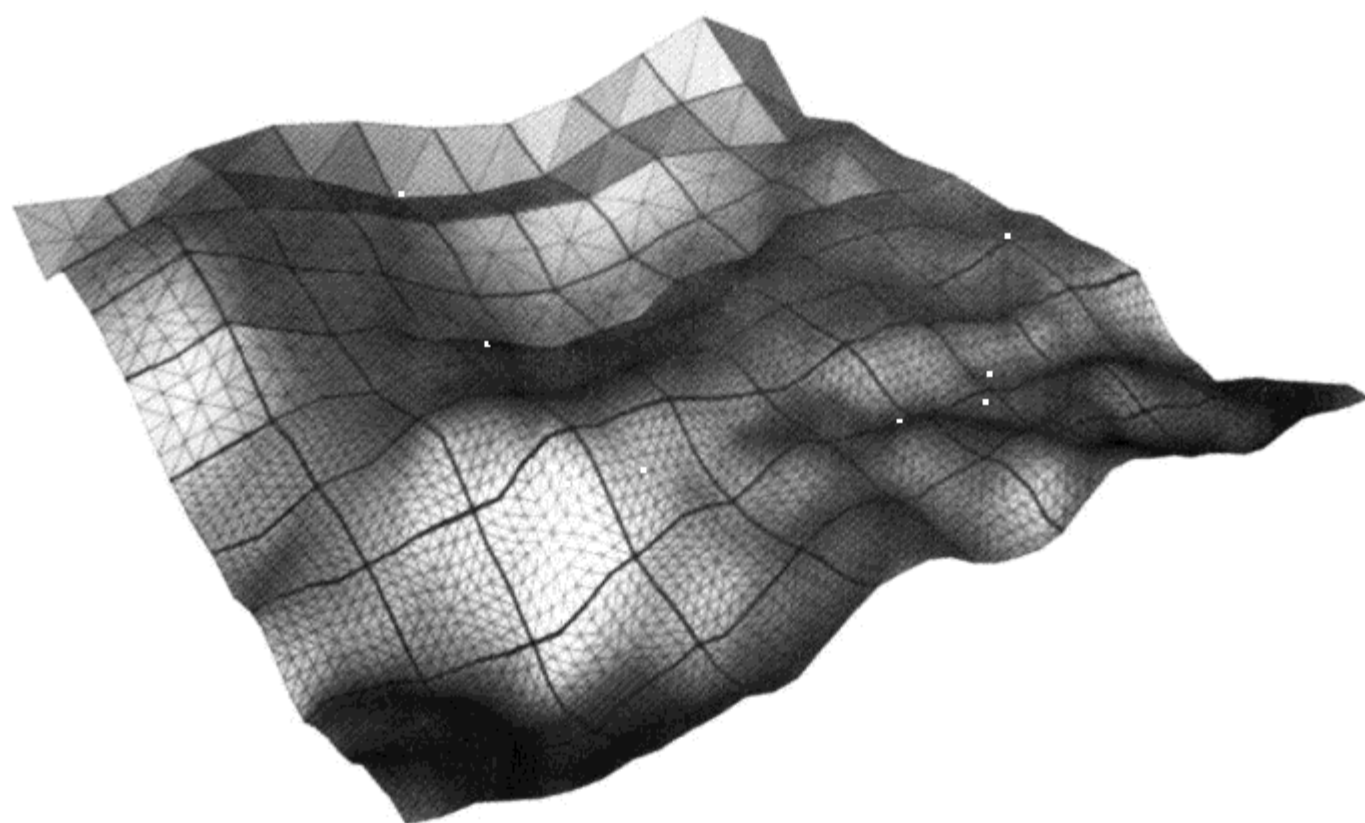
[Sarle01] Sarle, W. S., "Neural Network FAQ" available online at [www.ci.tuwien.ac.at/docs/services/nnfaq/FAQ.html](http://www.ci.tuwien.ac.at/docs/services/nnfaq/FAQ.html), February 15, 2001.







## 几何体管理





## 绪 论

Eric Lengyel, C4 Engine

lengyel@c4engine.com

这部分文章提供处理渲染到屏幕上的 3D 模型的几何体的技术。这个几何体可以在多种不同层次上进行操纵。本部分的一些文章描述执行于多边形和顶点层次以获取特效的技术。其他文章描述用于大型几何体管理和细节层次算法的方法。

可以认为,几何体管理是影响 3D 引擎性能最重要的因素。决定必须发送到图形卡的最小数据量(既对于可见几何体,也对于细节层次),产生的速度差异比其他任何方面都大。Tom Forsyth 的技术文章《各种 VIPM 方法的比较》以及 Greg Snook 的技术文章《使用联锁分片(Interlocking Tile)简化地形》均描述了平滑地减小场景中细节层次的技术。

用于组织场景中几何体的结构与场景本身中的物体一样重要。Ratcliff 的技术文章《快速可视剔除、射线跟踪以及范围搜索的球形树》描述了在球形树(SphereTree)内控制移动物体的方法,这个方法可用于可见性测试。Miguel Gomez 的技术文章《压缩的轴向包围盒树》以及 Matt Pritchard 的技术文章《直接访问四叉树查表》详细描述了使组织结构更小、更快而且更优的技术。

顶点和多边形操纵技术能够实现无穷的特效和易用技巧。Alex Vlachos 和 Evan Hart 撰写了两篇技术文章《近似鱼缸折射》和《渲染打印分辨率的屏幕快照》,它们描述了如何修改模型视图和投影矩阵(projection matrix),基于特殊目的更改顶点位置。Eric Lengyel 的文章《对任意表面应用贴花》演示如何通过剪裁表面的多边形网格在模型上应用表面印痕。

这部分内容还包括几篇其他方面的技术文章。Jason Shankel 的文章《用天空包围盒渲染远景》描述了搭建包含远景的天空包围盒(skybox)的信息。Alex Vlachos、David Gosselin 和 Jason Mitchell 撰写的文章《自阴影角色》介绍了一种角色自我产生遮光效果的技术。最后,Steve Rabin 的技术文章《经典 Super Mario 64 游戏第三人称控制和动画》描述了用户输入可以控制动画角色几何体的方法。

资源分享  
PDG

## 4.1 各种 VIPM 方法的比较

Tom Forsyth, Mucky Foot Productions  
tomf@muckyfoot.com

**视**点独立渐进网格技术 (View-Independent Progressive Meshing, VIPM) 由一项吸引了广泛研究兴趣的课题, 发展成了一项具有广阔前景的新技术, 成为所有优秀引擎的常备技术。现在 VIPM 技术已经融入到 Direct3D 图形 API 之中, 它几乎已经成为所有引擎不可或缺的一部分。在 Direct3DX 库中包含 VIPM 后, 再添加某种形式的 VIPM 就相当容易了。

然而, 在改善 VIPM 性能的努力中, 尤其是在尽可能提高硬件效率的努力中, 人们设计出了很多新的 VIPM 形式, 每种形式均具有各自的利弊和特征。本篇技术精粹作为读者开发更高级版本的铺路石, 帮助大家根据具体情况选择不同的变化版本。

本文假定读者对 VIPM 已有基本了解, 因此这里对 VIPM 不再进行详细的介绍。不过, 此方面已经出版了若干本优秀指导书, 而且可以在线阅读。其中两篇最著名的是 Jan Svarovsky 在《游戏编程精粹 i》[Svarovsky00]和 Charles Bloom 的网站上[Bloom01]发表的技术文章, 它们都是实践“vanilla”VIPM 方法的优秀自学资料。这里讨论的所有方法使用的都是同一种基本的塌陷/分裂 (collapse/split) 算法, 只是具体的实现方法不同。

### 4.1.1 考虑因素

在使用不同的方法时, 要考虑几个要点。不同的情况需做出不同的选择, 针对游戏中各种对象类型的不同使用方法, 也要相应地采用不同的 VIPM 方法。主要有以下几点需要考虑:

- **全局内存消耗** 模型的网格表达需要消耗多大的内存? 该内存为屏幕上显示的所有实例共享。
- **实例内存消耗** 屏幕上显示的每一个对象实例需要占用多少内存? 这部分内存空间要给每个实例复制一份, 不能共享。
- **流或引用的内存消耗** 这是每一帧实际引用的数据总量。对于对象来说, 可能有大量的数据保留在 CD 上, 或通过虚拟内存交换到某个硬件驱动器中。但是每帧实际引用的数据量可能并不大, 因此可以让数据流出去, 或由虚拟内存处理。这对内存空间有限的控制台来说尤其重要。
- **CPU 消耗** 在用户代码方面, 该算法使用多少个时钟周期? 包括单帧渲染开销以及帧与帧之间过渡时改变细节层次的开销。

- **API 接口效率** 驱动器和 API 接口将数据载入到图形卡时需要多少时钟周期?
- **总线带宽** 有多少数据需要发送到图形卡? 在 PC 机上, 这就是 AGP 总线带宽。
- **顶点缓存相关性** 尽管每个顶点被多个三角形使用, 但现代图形卡力图对每个顶点

只进行一次获取、变换以及点亮操作。为了实现这一点, 我们设计了一个顶点缓存, 用于保存最近使用的顶点, 应用程序需要尽可能使用该缓存里的顶点, 以获得最佳性能。这样, 当算法具有较高的顶点缓存命中率时, 即便使用了较多的三角形, 算法速度可能仍然较快。

顶点缓存相关性可以用每绘制一个三角形需要载入和处理的顶点数 (或简称“每三角形的顶点数”) 来表征。现有的静态网格 (非 VIPM) 三角形重排算法使用大约 16 项的现代缓存, 这个量大致可以降到 0.65。参考文献[Hoppe99]提供了一个例子。可以以此为基准衡量网格转换为 VIPM 时的有效性。还需要注意的是, 当使用三角形条带计算每个三角形的顶点数时, 只能计算画出来的三角形, 已经退化的三角形不能计算入内。退化的三角形对场景不能做任何贡献。

擅长处理数据流的算法能让应用程序绘制出宏大的场景, 这些场景的数据主要存储在磁盘上。如果数据流在途中某处达到极限, 如磁盘带宽、机器可用内存容量等极限值, 此时, 算法会适度地降低图像质量。

这种方法对拥有虚拟内存的系统也有一定帮助。当线性访问数据时, 虚拟内存管理器能把已经访问, 或长时间没被访问的数据交换出来。它甚至能进一步优化静态数据, 使之成为只读内存映射文件。这也能保证讨厌的“载入级别”信息尽量减少。对象的有关数据不需要在一开始全部载入。开始的时候, 玩家可以仅使用一些低分辨率的数据, 随着对细节模型需要的增加, 数据再逐步载入。

本文讨论的所有方法都基于同一个基本算法的实现。仅通过单一的操作就可以沿着该顶点的三条边中的一条将一个顶点塌陷 (collapse) 到另一个顶点上去。不产生新的“平均”顶点, 也不允许在没有三角形的边相连接的顶点之间进行塌陷。以上是一些值得注意的地方。不过, 目前的一致看法是, 在当前大多数硬件环境下, 对于等价的错误级别, 以上算法将耗费较高的运行时开销。当然, 随着技术的不断发展, 事情将发生改变, 新算法也会层出不穷。

对下面所使用的术语作一些说明: 网格的分辨率和其中的三角形数量成正比。因此, 高分辨率网格经地过边塌陷, 变成低分辨率的网格。与边塌陷相反的操作是边分裂 (split)。在边分裂操作中, 一个顶点分裂成两个独立的顶点。对于一个已知边的塌陷, 一个顶点称作保持顶点 (kept vertex), 另一个称为收藏顶点 (binned vertex)。在低分辨率网格中不使用收藏顶点, 仅使用保持顶点。在边塌陷操作中, 也有两类三角形。含有被塌陷边的三角形不出现在任何低分辨率网格中, 即被收藏了。在典型的塌陷操作中, 一般会有两个已收藏三角形, 当然, 在复杂的网络拓扑结构, 收藏三角形的数目会略有不同。那些没有收藏但使用收藏顶点的三角形被称为“变换”三角形, 通过变换, 这些三角中的收藏顶点用保持顶点来代替。执行边分裂操作时, 先前收藏的顶点和三角形变成“新的”了, 不过, 它们仍然称做收藏顶点和三角形, 这是因为一般情况下不存在分裂数据结构, 只存在塌陷逆过程的数据结构。由于大多数透视都是沿着塌陷的方向, 所以, “首先”、“接着”、“之前”、“之后”等词汇都假定是针对从高分辨率的网格塌陷到低分辨率的网格的过程。同样, 通过塌陷逆操作来完成分裂。

本篇技术以 PC 和 DirectX 为中心, 讲解 CPU、AGP 总线、图形卡、系统/视频/AGP 存储器、索引以及顶点缓冲。这样安排主要是出于讲解上的方便, 因为大多数控制台都有等价

的单元和概念。在存在显著差异的地方，我们都会予以强调。读者可能对术语“AGP 总线”不太熟悉。AGP 总线连接主系统存储器（以及 CPU）和图形卡（包括它的存储器）。AGP 总线有多种速率，通常在 500MB/s 左右。和连接 CPU 与系统存储器的总线，以及连接图形芯片和它的视频存储器之间的总线速度相比，这个速度相当小。一些控制台也有类似的瓶颈，其他控制台使用统一的存储器规划来避免该瓶颈的出现。在很多情况下，AGP 总线的速度瓶颈是 PC 图形处理的限制因素。

#### 4.1.2 Vanilla VIPM

Vanilla VIPM 是一种最著名的 VIPM，该技术被 Direct3DX8 库采用。它包含静态顶点的全局列表，该列表按照“先是尾收藏顶点，再是首收藏顶点”的顺序排列。每进行一次塌陷，被收藏的顶点为表尾顶点，同时，所使用的顶点数减少 1。这样可以保证使用的顶点在一个连续块中总是处于顶点缓冲的起点，这意味着，线性软件 T&L 管道总是只处理使用中的顶点。

三角形也是按照“先是尾收藏三角形，再是首收藏三角形”的顺序排列。每一次边塌陷操作一般删除两个三角形，虽然这些三角形实际上可以在任何地方删除，从零直到具有复杂拓扑结构的网格。

在边塌陷操作中，只被变换而没有被收藏的那些三角形具有由收藏状态变换到保持状态的顶点的索引。由于索引表随着细节层次的变化而变化，三角形索引缓冲将作为每个实例的数据存储。索引缓冲是由三角形的索引表组成（每个三角形由三个独立的索引确定），而不是由编入索引的三角形条带组成。

每个塌陷数据的记录格式如下：

```
struct VanillaCollapseRecord
{
    // The offset of the vertex that doesn't vanish/appear.
    unsigned short wKeptVert;
    // Number of tris removed/added.
    unsigned char  bNumTris;
    // How many entries in wIndexOffset[].
    unsigned char  bNumChanges;
    // How many entries in wIndexOffset[] in the previous action.
    unsigned char  bPrevNumChanges;
    // Packing to get correct short alignment.
    unsigned char  bPadding[1];

    // The offsets of the indices to change.
    // This will be of actual length bNumChanges,
    // then immediately after in memory will be the next record.
    unsigned short wIndexOffset[];
};
```

这个结构并非定长，wIndexOffset[]可以增长到需要变化的顶点数量。这样的结构从一定程度上导致了访问函数的复杂化，但是保证在进行塌陷或分裂时，所有的塌陷数据按照连续存储器地址存放。这样的存储方式能提高缓存线和缓冲预读取算法的效率。也使得数据流动



或从磁盘上载入塌陷数据的操作变得更容易。由于这些数据是静态的全局数据，它们也可以做成只读的存储器映射文件，这在很多操作系统下将大大提高效率。

虽然，初看 `bPrevNumChanges` 之时，似乎塌陷操作并不需要它，但在执行分裂和退回列表的时候需要它。需要先前结构中 `wIndexOffset[]` 的项数，用于跳过这些项。虽然这样做会使 C 程序代码看起来相当复杂，但生成的汇编代码实际上非常简单。

在进行塌陷操作时，由于被收藏的顶点总是处于列表的尾部，所以所使用的顶点数是减少的。三角形的数量减少 `bNumTris` 个，同样，收藏三角形也总是位于列表的尾部。

为了用保持顶点 (`kept vertex`) 代替收藏顶点，所有变换的三角形都需要重定向，对应于收藏顶点的索引变化保存在 `wIndexOffset[]` 中。每一个 `wIndexOffset[]` 都对应于一个需要变换的索引，它从收藏顶点的索引 (总是为最后一个) 变换成保持顶点的索引 (`wKeptVert`)。

```
VanillaCollapseRecord *pVCRCur = the current collapse;
iCurNumVerts--;
iCurNumTris -= pVCRCur->bNumTris;

unsigned short *pwIndices;
// Get the pointer to the instance index buffer.
pIndexBuffer->Lock ( &pwIndices );
for ( int i = 0; i < pVCRCur->bNumChanges; i++ )
{
    ASSERT ( pwIndices[pVCRCur->wIndexOffset[i]] ==
              (unsigned short)iCurNumVerts );
    pwIndices[pVCRCur->wIndexOffset[i]] = pVCRCur->wKeptVert;
}
// Give the index buffer back to the hardware.
pIndexBuffer->Unlock();
// Remember, it's not a simple ++
// (though the operator could be overloaded).
pVCRCur = pVCRCur->Next();
```

应注意，在某些体系结构中，从硬件索引缓冲读数据是一个不太理想的操作，所以要格外小心 `ASSERT()` 在干什么——它主要是为了演示之用 (如图 4.1.1 所示)。

分裂是塌陷的逆操作：

```
VanillaCollapseRecord *pVCRCur = the current collapse;
pVCRCur = pVCRCur->Prev();
unsigned short *pwIndices;
pIndexBuffer->Lock ( &pwIndices );
for ( int i = 0; i < pVCRCur->bNumChanges; i++ )
{
    ASSERT ( pwIndices[pVCRCur->wIndexOffset[i]] ==
              pVCRCur->wKeptVert );
    pwIndices[pVCRCur->wIndexOffset[i]] =
        (unsigned short)iCurNumVerts;
}
iCurNumTris += pVCRCur->bNumTris;
iCurNumVerts++;
```

```
pIndexBuffer->Unlock();
```

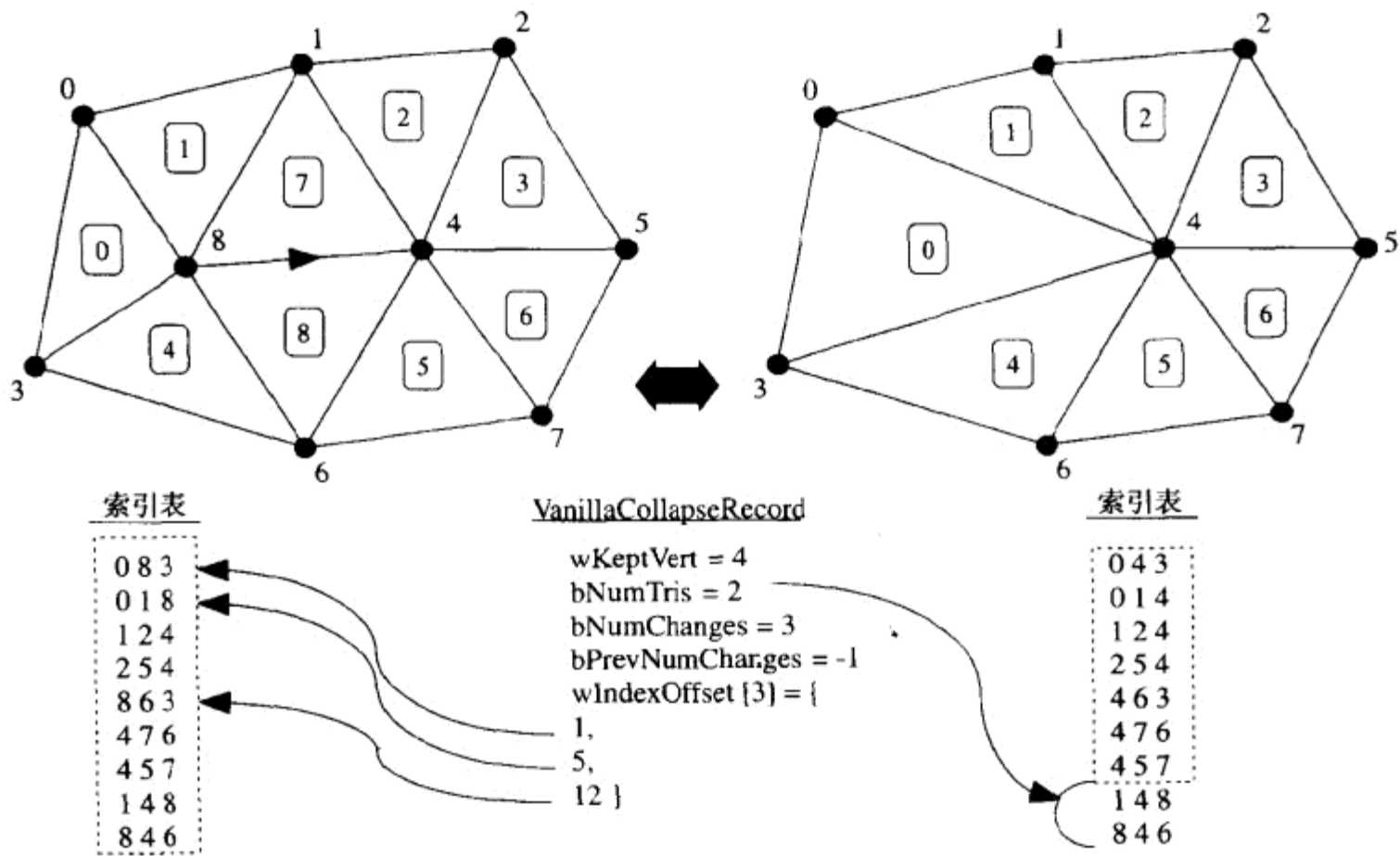


图 4.1.1 边的塌陷及塌陷前后索引表及 VanillaCollapseRecord

请注意，在实际应用中，出于某些历史原因，在示例代码里，VertexCollapseRecords 是从后往前存的，所以要交换 Prev() 和 Next()的调用。

vanilla VIPM 简单、易于编程，运行速度也相当不错。它大概是评价 VIPM 所使用的第一个版本，因为它是如此的简单，甚至还具有相当好的可扩展性、数据可流动性等等。

vanilla VIPM 的一个优点是，它能很好的实现数据的流动。塌陷信息和索引缓冲数据在存储器中完全呈线性，而且按塌陷排序。因此，当不是立即使用数据时，实现具有回退功能的数据流系统就极其容易。

然而，vanilla VIPM 技术也有一些不足。顶点缓存的相关性较差。由于三角形顺序是严格按照塌陷顺序决定的，所以，无法重排三角形以获取更好的顶点缓存性能。

另一个问题是，每个实例占用的存储空间相当大。对每一个实例，都需要复制整个索引数据块。只分配和实际应用数目一样多的索引，并且按需要进行伸缩，可以减少内存消耗（同时避免总是一直调用 malloc()和 free()），但是如果屏幕上显示大量对象时，内存的消耗量仍然很大。

最后，vanilla VIPM 只能和索引三角形表一起工作，这对适用条带数据结构的硬件系统来说是一个不够理想的选择。

### 4.1.3 跳带

跳带 (skip strip) 是一个带有一些重载性质的术语。它来源于与视点依赖渐进网格 (View-Dependent Progressive Meshing, VDPM) 相关的一篇论文[El-Sana99]。VDPM 要复杂得多, 而且要获得高效, 需要相当广泛的数据结构, 跳表 (skip list) 就是其中的一种数据结构。不过, 正是由于注意到收藏三角形时, 不必将其像 vanilla VIPM 一样从索引表的末端去掉, 因此产生了跳带技术。然而, 如果仅仅将一个顶点 (通常是收藏顶点) 移动到另一个顶点 (通常是保持顶点), 并把它留在已绘三角形列表中, 来使之退化, 这样就不会发生太多错误。硬件能很快定位退化三角形, 不必绘制任何像素, 就可以将它们剔除。

这意味着, 三角形的顺序不再遵循塌陷次序, 而且可以依照其他标准来排序。最初的跳带论文指出的一个精妙之处是, 现在三角形可以按照带序 (strip order) 排序, 而且实际上可以转换成为带。适合对数据按带序组织的硬件来说, 这是一个优势。由于这种 VIPM 方法是根据这篇论文得来的, 所以也继承了这个名字, 尽管有点不贴切。

重排三角形的能力提高了顶点缓存的相关性。条带自然适宜于此, 绘制一个三角形平均引用 1.0 个顶点 (对没有退化的长条带来说)。如果能选择合适的排序以及具有合适大小的顶点缓存, 这个值还能低得多。

实现代码的一个精妙之处在于, 塌陷/分裂程序和数据结构与 vanilla VIPM 几乎完全一样。惟一的不同点的是, 绘制三角形的数量并不随塌陷或分裂而改变。三角形只是简单地退化, 而不是从列表的末端删除。

然而, 这也显示了跳带存在的一个大问题。经过多次塌陷以后, 在列表中会有很多退化了的三角形。虽然它们能很快被硬件过滤, 但过滤的过程仍需要耗费一定的时间, 而且它们的索引数据同样要传送到图形卡中。这会占用总线带宽, 从而减少每秒可见的三角形数量。

经过多次塌陷以后, 顶点缓存的效率也会下降。干净整洁的跳带会因塌陷操作而弯曲、破坏, 从而导致缓存效率的降低。此外, 随着三角形的退化, 指向剩余顶点的索引数会增加。而收藏该顶点的塌陷必须更改引用该顶点 (包括退化三角形) 的所有索引。因此, 由于 `wIndexOffset[]` 的尺寸变大了, 进行的塌陷操作越多, 每次塌陷操作的开销就会越大, 操作的开销并不和绘制三角形的数目成比例, 这是不愿意看到的, 因为 VIPM 技术的出发点就是让较低细节的对象只花费较少的时间去渲染。

### 4.1.4 多层跳带

幸运的是, 跳跃条带结构的以上缺陷并不是不能解决的。经过一定次数的塌陷操作后, 暂停处理过程, 取得完成了所有塌陷操作的几何图体, 删除所有退化三角形, 并重新生成一个全新的跳带。然后通过此跳带继续塌陷操作, 直至其效率下降到极限, 如此等等。

每次创建一个新跳带层时, 都要将退化三角形全部删除, 这样可以降低传送到图形卡的三角形数量 (包括可见的和不可见的)。此外, 还要重新组织三角形列表, 以实现顶点缓存的优化。新塌陷操作不需要每次大量改变退化三角形的索引, 每个实例只需要拷贝实际使用的跳带层, 而且, 随着细节的降低, 它们逐渐变短。

不同的索引表可以全局存储，因为每次转换到一个新的列表，都会产生一个新的拷贝，并通过塌陷精简到所需三角形的数量。因此，此时多索引表的存在不再是累赘——它们是全局数据。这样也能恢复一些 vanilla 方法具有的良好数据流友好性。尽管图形的粒度 (levels of granularity) 会变得稍微粗些，必须在依照某层次渲染前获取整个索引表，但至少这不再是或全取或不取的事情，而且较低分辨率的索引表实际上很小。

为了稍微提高效率，完全塌陷（在转换成更低分辨率表之前）和完全不塌陷这两个版本的索引表都可以存储在全局空间里。这意味着，在两个索引表的边界两边摆动的单个塌陷仍然相当有效。如果只使用不塌陷的版本，每当细节层次增加时，就必须拷贝更高分辨率的索引表，然后，必须执行所有塌陷来画出下一帧。而存储塌陷版本意味着， $n$  次塌陷细节层次里的变化实际上只需要  $n$  次塌陷（有时更少）。

除了用一个全局结构的数组保存预先做好的索引表、对应每个索引表的塌陷表以及其中的塌陷数之外，实际塌陷/分裂的代码和结构都与标准跳带一样。在做任何塌陷和分裂操作之前，程序检查是否需要变化层次，如果需要，就拷贝新层次索引表，并开始执行塌陷/分裂操作，直至达到正确的细节层次。

这样一来，和 vanilla 比起来，通过增加全局存储（但易于数据流动或交换），以上技术弥补了跳带的所有缺陷。

跳带也能使用三角形表代替三角形条带，二者原理是完全一样的，但是使用的基本元素不一样。一些算法要求用表而不是条带，而一些顶点缓存程序用表代替条带可以获得更高的缓存速率。在给出的代码举例里，我们没有提供单独的实现，因为二者太相似了。

#### 4.1.5 混合模式 VIPM

到目前为止所提到的 VIPM 类型有一个问题，即需要复制对象的每一个实例的整个索引表。在某些情况下，这可能会造成相当大的负担，特别是在内存有限、控制台特别的机器之上，硬要把所有东西都塞进资源相对十分紧张的内存中时，有时甚至是在 VIPM 被引用之前，这的确会给系统带来很大的负担。如果将这些索引表的某些部分移到全局内存（即实例间的静态和共享内存），而不必为每个实例做一份拷贝，情况就会好得多。

在多级跳带中，即使那一级完全被塌陷，也会有很多三角形并没有受到影响。因此，不需要逐个实例地拷贝这些三角形；他们可以是全局，并且在实例之间共享。实际上，对于这个算法，我们使用的是索引表，有关索引条的情况会在稍后作为一个变体来讨论。在每一个层次上，三角形都分为如下 4 种列表：

- 没有被塌陷影响的三角形；
- 被塌陷操作收藏，但是在这之前没有被任何操作修改的三角形；
- 被塌陷修改，但是没有被收藏的三角形；
- 先被一次和多次塌陷修改，然后被收藏的三角形。

列表 2 和 4 都按收藏顺序排序，这一点和 vanilla VIPM 是一致的。列表 1 和 3 依照能产生最高顶点缓存效率的顺序排序。列表 2 附加在列表 1 之后，组合的列表放入一个静态、被所有实例共享的全局索引缓冲中。列表 4 则附加在列表 3 之后，在使用该层次的时候，将组合动态表拷贝到实例里面。这个列表在运行时，按照和 vanilla VIPM 完全一样的修改算法进



行修改。

为了绘制网格，对动态的每个实例的表，完成必需的塌陷和分裂，并绘制该表。然后，绘制相关联层次的静态表，惟一的修改之处是，塌陷静态三角形时，绘制三角形的数目有所改变。

所需的代码和结构基于多层跳表，只是每个层有两个表：一个拷贝的动态表，一个共享的静态表。另一个变化是，有两个三角形计数，每个表一个，而且一次塌陷可以同时改变两个计数，或者改变其中的任何一个。因此，成员 `bNumTris` 被 `bNumStaticTris` 和 `bNumDynamicTris` 代替，同时还要添加适当的增量和减量。

这意味着，每个网格的很大一部分是从一个针对缓冲相关性优化过的静态索引缓冲（表1）绘制而来的。它并不像想像中的那么好，因为这个表里的三角形只是构造该对象的一部分。在网格中三角形被移到其他3个表后会出现很多“空洞”，这样不仅减少最大顶点数量，也减少了每个三角形的实际顶点数量。尽管对某些动态缓冲排序来优化顶点缓存行为（表3），但由于塌陷会影响优化效果，而且表3的网格一般来说没有实现有效连接，所以对任何重新排序来说，能起的作用是有限的。

正如所有的多层方法那样，该方法的数据流是友好的。然而在这种情况下，由于表是按照塌陷顺序组织的，粒度在三角形层次变得更精细了，而不仅仅在表的层次。这个结果是否十分令人激动是另外一个问题——因为更精细的控制很可能并没有在性能上产生很大的差异。

这确实需要两次调用 `Direct3D`（或者等效的 `API`）的 `DrawIndexedPrimitive()`。虽然在大多数平台下，这不是一个瓶颈，而且并不影响渲染速度。但是对非常少的三角形网格来说就很重要，对这些网格，也许适合于转换为另外一种方法。

#### 4.1.6 混合模式跳带

除了使用的是条带，而不是用 `vanilla VIPM` 的动态表实现等差异以外，混合模式跳带和混合模式表是完全一样的。对于跳带来说，使用条带意味着按照塌陷顺序排序效率太低，也就是说，表2的三角形现在必须通过退化而被收藏。这就迫使这些三角形变成动态的而不是静态，而且他们加入到表3和表4。这3个表中的三角形被合并，并作为一个跳带处理——重排获得最优缓存效率，拷贝每个实例，按照塌陷信息进行修正。

这个方法的缺点是，对于每个实例来说，现在需要拷贝的数据更多了，而且因为三角形是以带序排列的，而不是按塌陷顺序排列，所以，简单地把他们从索引表尾部删掉，并不能完全收藏三角形。不过，这些不利因素只是比使用表的版本略微差了一点，如果硬件需求适合条带结构，这仍然是一个很不错的方法。

#### 4.1.7 滑窗

滑窗 `VIPM` 引入了完全静态和全局的索引缓冲思想，不需要编辑索引，因此每个实例只需要很小的内存空间。

滑窗算法注意到，当发生塌陷时，会有两种类型的三角形：收藏三角形和修正三角形。

不过，对于修正三角形来说，在塌陷前和塌陷后在索引缓冲的同一个物理位置上出现是没有实际必要的。三角形的旧形式可以与收藏三角形一起，简单地从索引缓冲的尾端丢弃，而新版本添加到另一端。

因此，在例子中，塌陷操作不是收藏两个三角形、编辑其他 3 个三角形，而实际上收藏 5 个三角形，并且加上 3 个新三角形。这两个操作都仅仅是通过更改用于渲染的首索引和尾索引，即让“渲染窗”沿着索引缓冲滑行（如图 4.1.2 所示）。

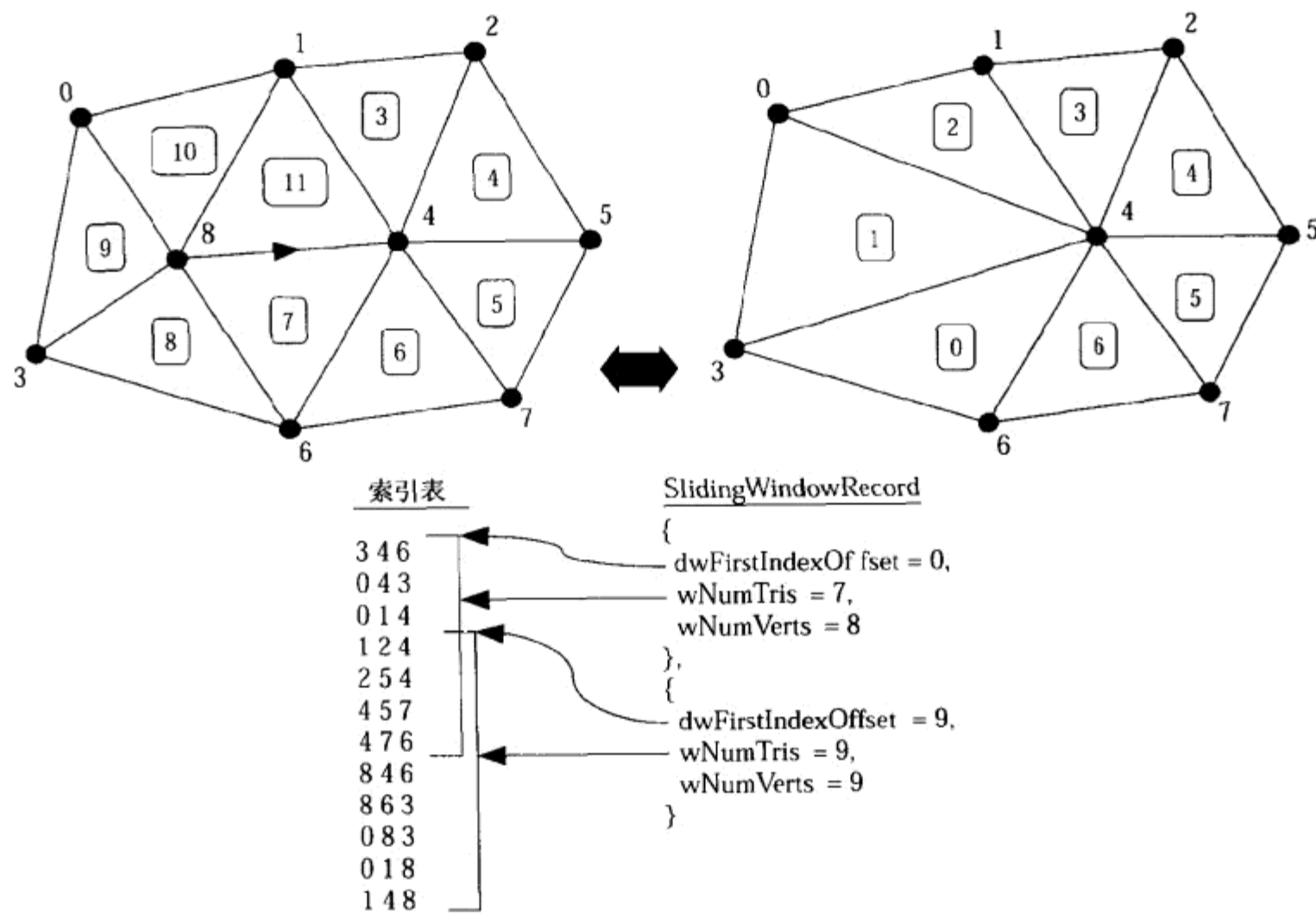


图 4.1.2 显示索引表和两个窗口的塌陷

索引缓冲区被分割为 3 个部分。最开始是变化所得的三角形，按照塌陷的反序排列。中间是没有被塌陷影响的三角形，以任意（顶点缓存优化）顺序排列。最后是在塌陷中被转变或者被收藏的三角形，同样以塌陷的反序排列，即第一次塌陷排列在末端。注意，作为塌陷结果而被修改的三角形不能再一次塌陷（不管是收藏还是变换）。如果一个三角形被第二次塌陷修改，就意味着这个三角形必须从索引缓冲的末端删除。但是，它已经被加到索引缓冲的开始，所以不能又从末端删除——通过排序来实现的机会极其小。

三角形一旦被塌陷操作修改，它可以被卷入另一次塌陷的惟一途径是，启动一个新索引缓冲，而且它和先前塌陷的索引缓冲正好有一样的三角形。这个新索引缓冲的顺序不受前一次塌陷的限制，所以可以通过新的塌陷来排序。这里再次用到了多层（multilevel）的概念，但是在这个情况下，如果不借助多层，就无法进一步塌陷，而不仅是出于效率方面的考虑。

这个问题的表面价值是，诸如 QEM 这样的算法为塌陷提供了一个排序。如果严格遵循这个顺序，QEM 经常需要执行新塌陷，对先前塌陷已经修改的三角形再做一次塌陷。这样强



迫生成一个新层，而且需要拷贝索引缓冲。因为只执行了很少几次塌陷，因此，这个拷贝几乎和原来的一样大。如果在产生一个拷贝之前，只执行了少数几次塌陷操作，所有这些索引缓冲所使用的内存将是巨大的。

然而，实际上并不需要严格遵循 QEM 决定的塌陷顺序。渐进网格不是一门严格的科学，因为除了照相机到物体的距离外，它忽视了其他任何影响因素，而且它的全部出发点只是充分满足视觉需要。因此，实际上无需严格遵循 QEM 决定的塌陷顺序可以进行一定程度的变通。

具体的实现方法是，遵循 QEM 塌陷的顺序，直到 QEM 决定要执行一次涉及已经被修改三角形的塌陷操作为止。执行这样的塌陷操作会强行增加一个新层，所以这样的操作应尽可能地拖延。暂时忽略这次塌陷操作，寻找一个不需要创建新层的最佳塌陷操作。比较这两种塌陷的误差，如果它们在一定的容许值内，不以严格的顺序执行也不会严重影响视觉质量，那么就执行那个不需要产生新层的塌陷操作。

一旦误差太大，那么首先执行错误的塌陷，就会严重影响图像的质量，算法有效性受到严重影响，并且创造了一个新层。现在，在拷贝发生之前，执行适当数量的塌陷，三角形的数量大大减少了，因此在塌陷到最小的细节层次之前，只需要很少的层次。

示例代码使用一个相当小的容许误差级，只有平均塌陷误差的 10%，不过即使是这么小的容许误差，也显著减少了层数。使用更大的容许误差，可以更多地减少对内存空间的需求，不过，只是针对某个时刻来讲情况如此。过一段时间后，算法就会运行到塌陷没有涉及到的三角形。在这个状况出现之前，大多数的网格只能损失大概 20% 的三角形，但内存使用仍然维持在切合实际的级别上。

因为索引和顶点表并没有运行时 (runtime) 的修改，所有的数据都可以是全局，所以每个实例的内存使用几乎为零。转换细节层次的 CPU 时间也几乎为零——每次用一个简单的表查找来决定使用的索引表，从该索引表取出的首索引和尾索引，还要决定使用多少个顶点。在实际应用中，索引表是连接在一起的，所以首索引也就暗示了使用的索引表。这个查找的表由如下结构组成：

```
struct SlidingWindowRecord
{
    unsigned int    dwFirstIndexOffset;
    unsigned short wNumTris;
    unsigned short wNumVerts;
};
```

虽然已知三角形和顶点的数量少于 64k (这是当前所有已知硬件的极限)，又由于索引表是很多表的串联，它的长度很容易超过 64k，所以需要 32 位。这确实意味着该结构以 8 字节很好地填充对齐。其渲染代码却相当简单：

```
SlidingWindowRecord &pswr = swrRecords[iLoD];
d3ddevice->DrawIndexedPrimitive (
    D3DPT_TRIANGLELIST, // Primitive type
    0,                  // First used vertex
    pswr->wNumVerts,     // Number of used vertices
    pswr->dwFirstIndexOffset, // First index
```

```
pswr->wNumTris ); // Number of triangles
```

和其他方法不同，这里并没有执行分裂和塌陷的代码，当前的细节层次仅仅是在每次渲染物体时查找 `SlidingWindowRecord` 表。这也意味着，通过硬件转换和照明卡，渲染物体所需的 CPU 时间是固定的，无论什么细节层次，渲染每个物体的 CPU 时间是个常数。在任何算法中，“恒定时间”一词总是意味着好的性能。

滑窗 VIPM 的主要问题是，在每层索引表的起始和末尾处，它强迫三角形的有序排列。这会产生两个效果：一会使得条带很难使用——只有三角形表真正能很好地处理固定的排序。二会影响顶点缓存的效率。

幸运的是，事情并不像看起来的那样糟。当执行边塌陷的时候，删除了所有使用收藏顶点的三角形，所以这些三角形都到三角形表的末端去了。一般这样的三角形有 5 到 7 个不等，它们在收藏顶点的周围组成一个三角形扇面。然后加上新形式的三角形。这些三角形需要一起载入到索引表的头部，一般是 3~5 个三角形，它们围绕保持顶点组成一个三角形扇面。这些扇面可以排序，以取得最好的缓存器相关性。由于索引表的中部未受影响，即未排序，所以可以用来对顶点缓存重排序。这样就能获得比 `vanilla` 更好的缓存相关性。虽然它距离理论上的理想情况仍然相当远，但是在实际应用中是可以接受的。

在每层使用更大的中间索引表（让每层的塌陷操作更少），可提升顶点缓存相关性。这需要花费更多的内存，不过由于有额外的性能提高，这样做是值得的，特别是当它是全局内存的时候。

需要条带而不是列表的硬件也仍然可以使用这个方法，虽然它需要退化很多三角形以合并不同的部分。在实际应用中，这个方法并不会增加所需的索引数量，相反，实际上减少了索引数——对于条，每个三角形一个索引，而表是三个。每个画出的三角形的顶点缓存效率完全一样。原三角形的吞吐量增加了很多（差不多两倍），但是因为所有这些额外的三角形都为退化三角形，大多数硬件会很快的剔除它们。如果要进行选择，到底使用哪种基本方法，取决于硬件是受索引带宽的限制（在这种情况下，条带是最优的），还是受三角形吞吐量的限制（在这种情况下，表是最优的）。

#### 4.1.8 小结

VIPM 似乎已经发展成熟。现在已经是主流技术，并已融入主要 API 中。对于离散物体，通过 CPU 补偿，在大多数获取视觉冲击的效果上，VIPM 超越了 VDPM 和静态细节层次的方法（不过，需要指出的是，VDPM 在大场景，特别是规则高度场（`regular-height-field`）中仍然占据着重要的地位）。此外，现在可供选择的方法很多，它们各有优缺点。创新当然不会到此为止——现在已经有一些作为未来研究的有趣途径了。不过，当选择实现哪个 VIPM 方法时，这篇文章能给很多问题和决策提供相当好的指导。

表 4.1.1 说明了每种方法的相对长处和弱点。注意，其中的“跳带”指多层跳带。单层跳带在实践中不切合实际，其原因参见前文。

表 4.1.1 各种 VIPM 方法的优缺点总结

	Vanilla	跳 带	混合模式	滑 窗
顶点缓存的利用	不良	优	良	良
全局内存的利用	低	中	中	高
实例内存的利用	高	高	中	低
改变细节层次 (LoD) 的 CPU 开销	中	中	中	微小
API 效率	良	良	良	优
表效率	不良	优	良	良

4.1.9 参考文献

[Svarovsky00] Svarovsky, Jan, "View-Independent Progressive Meshing," *Game Programming Gems*, Charles River Media, 2000, pp. 454~464.

[Bloom01] Bloom, Charles, VIPM tutorial, and various VIPM thoughts gathered from many sources. [www.cbloom.com/3d/index.html](http://www.cbloom.com/3d/index.html).

[Hoppe99] Hoppe, Hugues, "Optimization of Mesh Locality for Transparent Vertex Caching," *Computer Graphics (SIGGRAPH 1999 proceedings)* pp. 269~276. See also [www.research.microsoft.com/~hoppe/](http://www.research.microsoft.com/~hoppe/).

[El-Sana99] J. El-Sana, F. Evans, A. Varshney, S. Skiena, E. Azanli, "Efficiently Computing and Updating Triangle Strips for View-Dependent Rendering," *The Journal of Computer Aided Design*, vol. 32, no. 13, pp. 753~772. See also [www.cs.bgu.ac.il/~el-sana/publication.html](http://www.cs.bgu.ac.il/~el-sana/publication.html).



## 4.2 使用联锁分片简化地形

Greg Snook

gregsnook@home.com

随着近来 3D 渲染硬件的发展，似乎每个人开发的游戏都置身在宏伟的外景中。遥远的地平线以及连绵起伏的地形，这些画面在以前模糊不清，而且经过了剪辑，但现在却可以获得真实的效果了。游戏程序员一度使用的 BSP 树和基于范围的渲染方法现在已成昨日黄花，取而代之的是闪耀的 ROAM 和 VDPM 等技术。

实时最优适配网格 (Real-time Optimally Adapting Meshes, ROAM) [Duchaineau] 和视点依赖渐进网格 (View Dependent Progressive Meshes, VDPM) [Hoppe98] 在这些参考文献中已有详细描述，这里只对它们作简要介绍。对于那些不需要大量几何体的地形部分，如适当的平原或远处范围，这两种方法能适当地减少多边形的数量，因此也减少了渲染的负荷。反过来，对于距离摄像机较近的场景或粗糙的表面，需要更多的多边形，这两种方法也可以让它们具有更多的细节。实质上，它们是获得同一个简单目标的两种复杂途径：在需要的时候使用较多的多边形，而不需要的时候使用较少的多边形。

对于某些应用，麻烦在于，ROAM 和 VDPM 这样的方法要依赖于从程序上产生几何体，来获得从低到高细节场面的平滑过渡。ROAM 使用三角形交集的二叉树从已知高程场来构建实际地形几何体。VDPM 使用粗糙网格表示低细节地形，在需要时运用一系列连续的顶点分裂进一步划分出地形多边形，以此获得相似的效果。在大部分情况下，这些连续的三角形剖分影响了硬件传输和光照的速度优势发挥，这些优势依赖于静态几何体获得最佳速度。

主要原因是，这些方法工作性能几乎太优秀。它们可以把地形剖析至多层次，精选出哪些保留，哪些消隐。随时间流逝，地形几何体可能会发生许多细微的变化，对发生变化的地形，整个方法又会作重新处理。所以限制划分几何的层次，从而处理大面积的静态几何体，能够维持硬件良好的工作性能，也能灵活地处理地形随时间变化的问题。

这篇技术文章讲解的是一种更简单的方法，大部分应用通过少量的编码就可以运用。它的目的不是产生 ROAM 或 VDPM 能够获得的视觉效果，而是通过动态适配细节层次和动画灵活性创建简单的地形。它完成这一工作的同时保持了完全适合于硬件传输和光照的数据系统。

### 4.2.1 分片的重访问

曾几何时，游戏程序员使用 2D 分片 (tile) 来表示游戏场景。只有一个简单原因：较少的原图更容易创建和控制。早期游戏不具备供给大量像素数据的内存，因此，分片较小的图片来创建较大场面的假像。这些小分片也更容易绘制和移动，所以平滑滚动的 2D 游戏将更容易从  $32 \times 32$  像素的小分片创建得来。

这里介绍的地形方法的基本工作原理相同，将地形分解成更小、可再用的分片。也具有相同的优点：少量数据易于处理，可以最优地绘制，而且能更有效地使用内存。明显的差异是，分片内的像素数据将不再处理。地形分片用索引缓冲来表达，它将地形的顶点链接在一起。

将 3D 分片看作从上向下投影到地表的格栅。每个格栅块表示地形系统中的一个分片。在表面上，看不出地形分片重复的迹像，仿佛地形是相当随机划分的几何图形。但在场景后台，分片并重用了大量数据。

以顶点和索引缓冲的形式对待每个地形分片。虽然每个分片可能包含一套独有的顶点数据，但用于绘制分片的索引缓冲可以非常频繁地重复使用。事实上，通过仔细规划顶点数据，可以创建一个分片索引缓冲有限集在整个地形中使用。

由于分片的细节有限，这给完成这一任务带来了便利。首先，每个分片一定包含相同数量的顶点，它们构成与相邻三角形的共用边。这些顶点表示该分片可能含有的最多细节层次。其次这些顶点在  $x$ 、 $y$  平面中分布为有规律的方格。使用  $z$  表示顶点的海拔高度。最后，按相同的顺序保存每个分片的顶点，让索引缓冲可以运用于任何分片上。如图 4.2.1 中的分片所示。在图 4.2.1 中，用一个  $17 \times 17$  的顶点分片图显示每个顶点的网格排列位置，对应每个顶点都有一个取自地形图的惟一高度值。

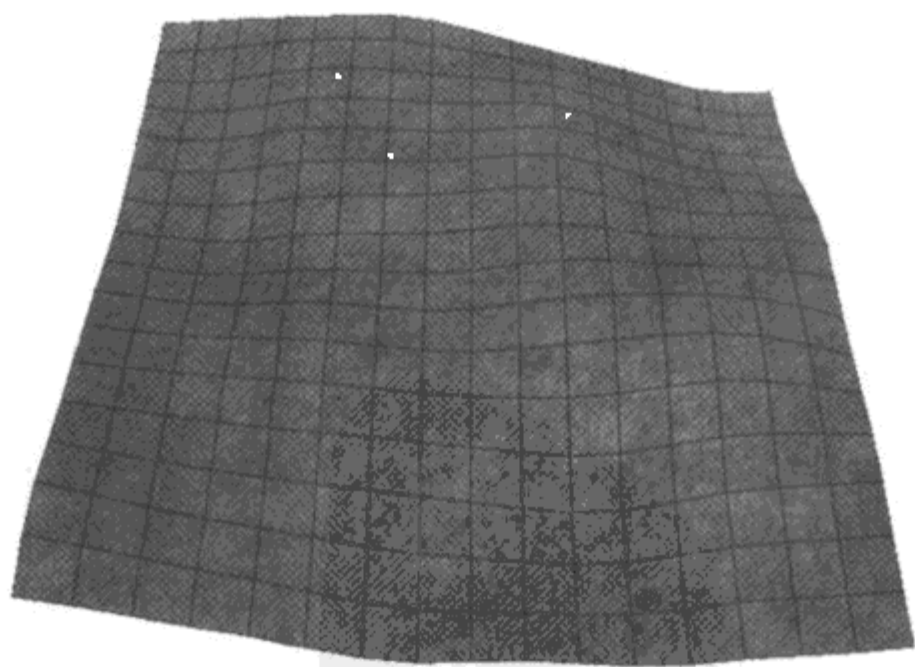


图 4.2.1 17×17 顶点的示例地形



采用这种顶点组织的原因很简单。因为顶点数据总是出现在规则网格内，而且顺序相同，对于整个地形可以创建固定的索引缓冲集。使用适当的索引缓冲，可以在任意层次上渲染给定的分片，从完全细节到简单的三角形对。使用更多顶点的索引缓冲表达更高细节的分片。相似地，使用较少顶点的索引缓冲以较少三角形渲染分片。图 4.2.2 通过不同细节层次的分片渲染演示了这一过程。

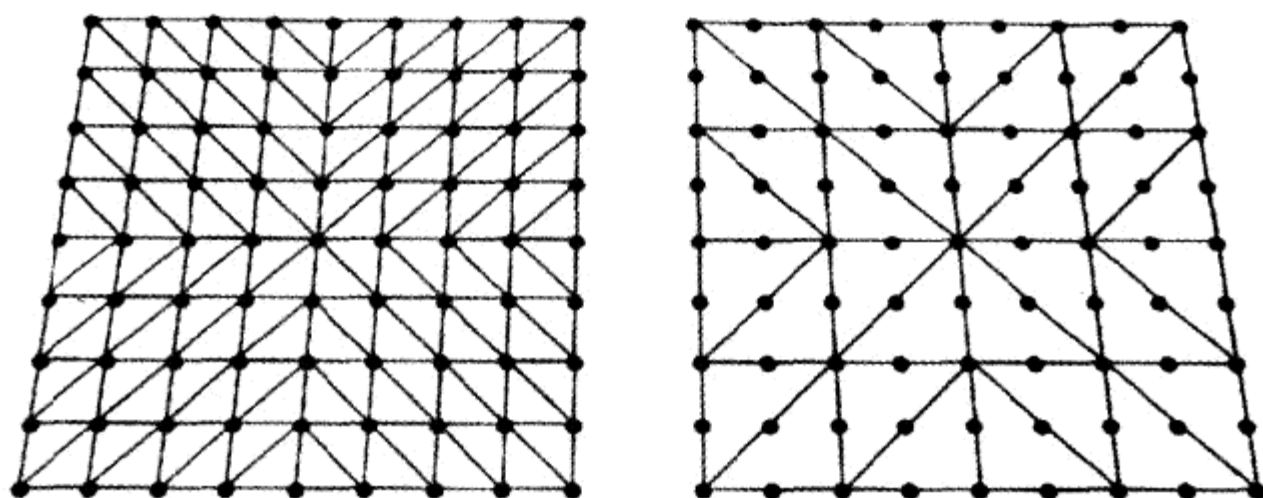


图 4.2.2 对相同顶点的分片使用索引缓冲来创建两个不同的细节层次

### 4.2.2 生成地图

为了创建地形分片，需要一套原始数据作为开始。一个常见方法是从高度图读取高度数据。这就是地形的灰度级位图，其中像素的亮度用来表示给定位置的高度。高度图另外的优点是它已经表示为有规律的网格。因此，它容易转换成地形顶点数据。它也可以作为动画资源，调整像素值来影响不同位置的地形高度。

创建分片顶点很简单。因为每个分片顶点在  $xy$  网格上有一个已知 2D 位置，剩下的就是从地形图取对应的高度像素，并转换成地形顶点的  $z$  值。对每个地形分片，可以取高度图的对应像素块，创建该分片的惟一顶点缓冲。在动画调试图情况下，这个过程必须周期性重复，以更新地形顶点。

### 4.2.3 分片模板

索引缓冲可以看作一个使用分片顶点的绘图模板。如图 4.2.2 所示，索引缓冲决定了如何从分片中取出三角形以控制绘制细节程度。根据主要规则，每个分片的顶点缓冲按相同的顺序布置，因此索引缓冲可以互换使用。对于一个  $9 \times 9$  顶点分片，可以创建一个全局索引缓冲集合来绘制  $9 \times 9$  顶点集合的所有可能细节层次，忽略网格中的顶点可创建不同级别的几何体。最高级的索引缓冲使用所有 81 个顶点来绘制 128 个三角形，而最低级索引缓冲只使用 4 个角顶点来绘制两个三角形的四边形。另外，还有两个分别表示 32 个和 8 个三角形的细节层次。





下一个需求是在 4 个细节层次中, 决定绘制每个分片时需要使用哪个细节层次。该决策可以来源于分片与摄像机之间路径的简单函数, 或者来源于对视角和觉察到的分片粗糙度的完全试探性考虑。所使用的最佳方法取决于游戏地形以及摄像机的移动。Hugues Hoppe 撰写的关于 VDPM 方法的文章[Hoppe]对用于为每个地形位置选择细节层次的试探性思路进行了更多的探讨。本书附带光盘上的应用示例 SimpleTerrain 使用分片到摄像机的距离进行简化。一旦确定了细节层次, 绘制工作就是一项将分片的顶点缓冲与期望的索引缓冲一起发送到选定的渲染 API 的简单工作。

#### 4.2.4 消除难看的接缝

现在已经具备了基本的地形系统, 但决不是一个平滑的地形。在现在拥有的地形中, 不同细节的分片之间存在突变。更有甚者, 在两个不同细节层次创建的分片之间可能出现裂缝。简短地说, 整个场景变得杂乱, 但可以想办法来补救。

方法的关键是拥有联锁的分片。也就是说, 创建完全啮合一起的分片, 尽管相邻分片之间的细节层次上有差异。为了达到此目的, 需要不同的索引缓冲集合将没有缝隙和接缝的不同层次分片合并在一起。这些索引缓冲可以分成两组: 主体和链接。主体表示已知细节层次分片的主要部分, 删除一些区域为链接片断提供空间。正如其名字所表明的, 链接将不同细节层次的主体无缝链接在一起。

图 4.2.3 显示了一个已知细节层次的 16 种可能主体类型。为了便于控制, 指定分片只向下链接, 即意味着更高层次的分片必须使用链接片断与较低细节层次的相邻分片连接在一起。如图 4.2.3 所示, 每种主体类型的无阴影区域表示的空间, 需要链接将较低细节层次相邻分片相连。

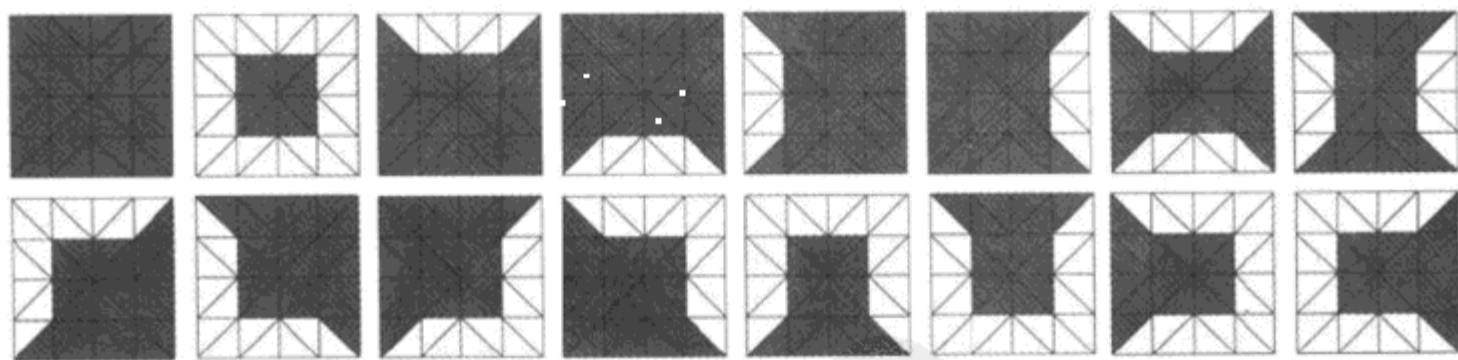


图 4.2.3 16 种基本分片主体。无阴影区域为必须放置链接片断的地方

链接片断为较小的索引缓冲, 它适应于主体分片留下的空白空间。这些索引缓冲从使用较多数量顶点的分片到使用较少顶点数量的相邻分片, 逐级安排三角形。图 4.2.4 是一个用于连接两个主体分片的示例链接。因为只向下链接细节层次, 所以每个细节层次需要足够的链接片断将它下面的细节层次连接在一起。例如, 对于  $9 \times 9$  的顶点分片, 对最高细节层次的每个边需要 3 个链接片, 因为它必须能够链接 3 个较低细节层次。最低细节层次 (即简单四边形分片) 不需要链接片, 因为所有的较高细节层次必须完成对它的链接。

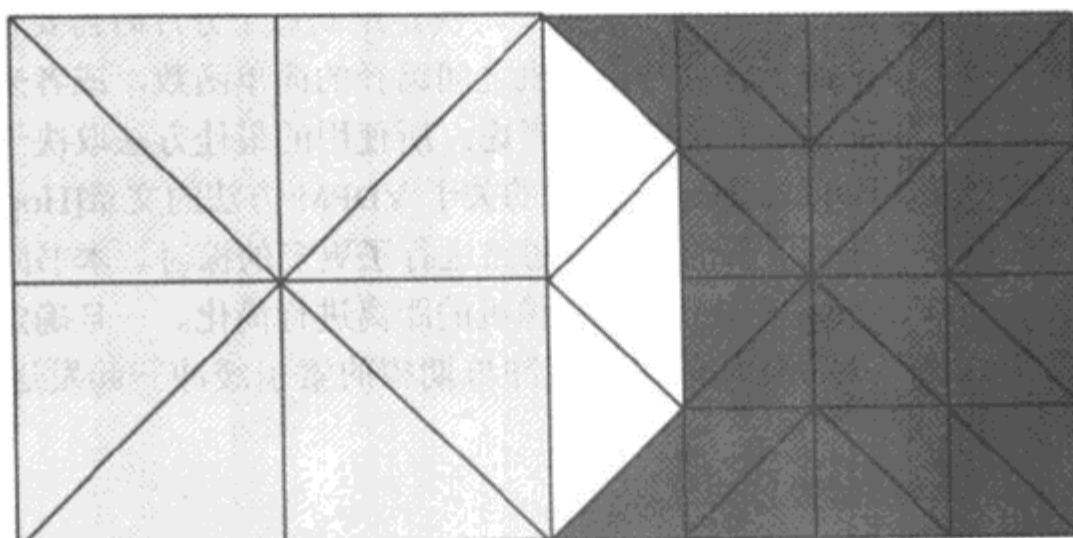


图 4.2.4 用于连接两个不同细节层次的链接片示例

已知  $9 \times 9$  顶点分片有 4 个细节层次，可以计算出所需的索引缓冲总数等于 48 个主体片与 24 个链接片的总和。表 4.2.1 显示了平滑地形所需的完整索引缓冲表。从表中可见，增加细节层次数也增加了索引缓冲数量，但因为它们的数量相当小，而且可以在整个地形中使用，因而仍然相当有效。

表 4.2.1 4 个细节层次的示例所需的所有索引缓存

细节层次	主体片断	+	链接片断	总和
4	16		3 (每边)	28
3	16		2 (每边)	24
2	15		1 (每边)	19
1	1		0	1
合计:				72

#### 4.2.5 更好、更快、更强

使用新主体和链接片，意味着需要更改渲染方法。对每个分片，现在需要检查分片的邻居。在每个与主体片相邻的低细节层次片一边，都保留缺口。然后，选择必需的链接片，填充这些缺口并与相邻分片链接在一起。将每一个这样的索引缓冲与该分片的顶点缓冲发送到渲染 API 进行绘制。在最坏情况下，每个分片发送 5 个索引缓冲（1 个主体、4 个链接），但在最佳情况下，只要发送 1 个索引缓冲（完全主体分片）。

将索引缓冲组织成三角形条带形和扇形，可以进一步优化这个方法。对于较大的分片（ $33 \times 33$  顶点以上），这将大大减少渲染时间。另外，可以调整分片中的顶点顺序获取渲染时的更佳性能。恰当的顺序取决于渲染分片时哪些索引缓冲用得最频繁。

### 4.2.6 结论



图 4.2.5 显示了渲染方法的最终输出。示例程序 SimpleTerrain 使用 DirectX 8.0 演示该方法。这个例程的完整源代码可以在附带光盘上找到。在这幅图中，显示了地形的栈框，以展示所使用的不同主体分片和链接分片。为了提高可读性，去掉了地面纹理。

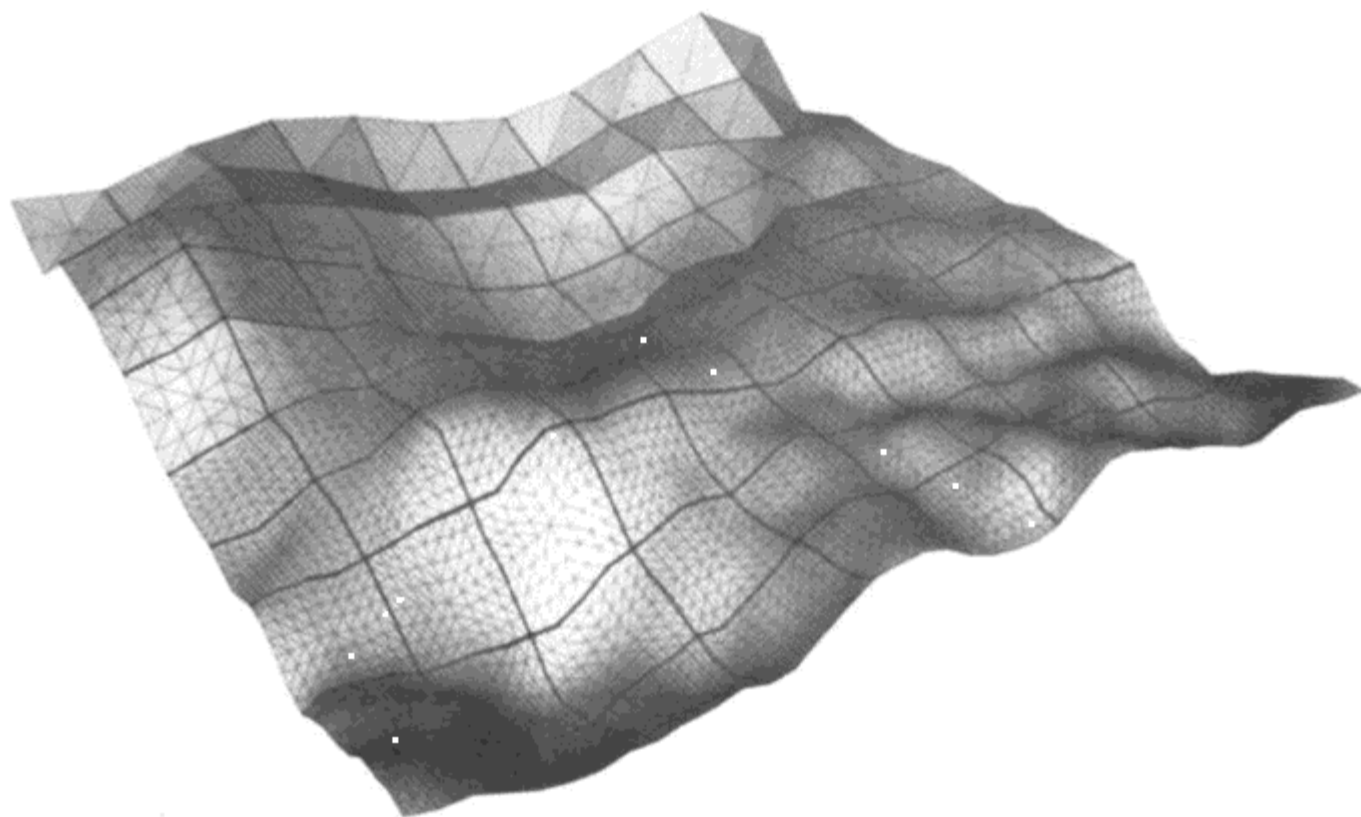


图 4.2.5 显示所使用分片的链接片断的 SimpleTerrain 示例程序的输出

这篇技术文章的目的是，为渲染动态地形提供一种替换流程序方法的方案，同时能完全实现硬件转换和光照。使用这个方法，可以实现动态地形系统，而且运行速度快，但不会严重影响应用的帧频。虽然最终的地形可能不像精心开发的 ROAM 或 VDPM 系统那样理想，但它提供了基本相同优点的同时，还可以提高硬件的渲染速度。

### 4.2.7 参考文献

[Duchaineau] Duchaineau, M., Wolinski, M., Sigeti, D., Miller, M., Aldrich, C., and Mineev-Weinstein, M., "ROAMing Terrain: Real-time Optimally Adapting Meshes" ([www.llnl.gov/graphics/ROAM](http://www.llnl.gov/graphics/ROAM)).

[Hoppe98] Hoppe, H. "Smooth View-Dependent Level-of-Detail Control and Its Application to Terrain Rendering" *IEEE Visualization 1998*, October 1998, pp. 35~42. ([www.research.microsoft.com/~hoppe](http://www.research.microsoft.com/~hoppe)).

## 4.3 快速可视剔除、射线跟踪以及范围搜索的球形树

---

John W. Ratcliff, Sony Online Entertainment  
jratcliff@verant.com

虽然保存静态 3D 物体有许多数据结构, 包括四叉树、八叉树以及 BSP 树, 但它们对于大量的动态物体并非总是理想的数据结构。这篇文章介绍一个算法, 并演示其控制上千个运动物体的应用, 这些物体连续不断地在球形树 (sphere tree) 被维护成层次包围球。

这个算法的设计目标是在 99% 的情况下, 在树结构中更新物体的移动所耗的 CPU 时间几乎为 0。虽然球形树中的执行的检测比其他数据结构多, 但维护树使用非常少的 CPU 时间, 因此可以得到补偿。这个数据结构非常适合于在庞大世界空间中对巨大数量移动物体进行剔除。对象以完全不同的速度移动, 甚至许多物体根本不移动, 也不会对它产生影响。当频繁地从树中插入和删除对象时, 它的开销也非常低。

### 4.3.1 包围球

---

当使用包围球 (bounding sphere) 作为一种剔除的基本方式时, 就会存在一些局限。包围球并不能非常紧密地包围一个物体, 尤其当物体又高又瘦的时候。另一方面, 以上描述的包围球可以看作一个特征, 不必是一种限制。包围球必须包含对象各个方向上的内容。这要考虑所有可能出现的动物动作。另外, 也假定这个包围球包含所有子物体。因此, 也可以假定父体的所有可见性状态对它的子体也成立。上述包围球的另一个优点是, 它可以用于剔除动画、阴影和其他相关效果。物体周围的超出范围可以是处理物体、它的子体以及相关效果的有效工具决定他们何时激活或静止。剔除阴影、动画以及特效象剔除几何体一样严格。

### 4.3.2 使用球形树

---

模拟中的每个物体, 不管它是运动的还是静止的, 都使用类 SpherePack 将其包在一个有效的球形树中。当物体使用方法 NewPos() 改变位置时, SpherePack 就计算新位置与父节点中心的平方距离。如果它仍然包含在父体球的半径内 (几乎所有时间均设计为真), 则程序立即返回。这个方法按内联实现, 以获得最佳性能。对于大多数运动的物体, 这是惟一执行的计算, 即使这些对象有上千个, 情况也是如此。对于静态物体, 除了初始插

入树中外，不执行任何计算。

当新位置导致物体穿过它的父体球表层时，将子体从父体中删除，并插入到树的根节点。这涉及到两个指针交换来即时维护完全有效的树。当子节点从父节点分离后，它将放到重组的 FIFO 队列（queue）中。将父节点添加到重新计算的 FIFO 队列中，维护一棵最优平衡树。在每一帧，模拟在 SpherePackFactory 中的执行处理方法，以至于可以刷新重组和重算的 FIFO 队列。

四叉树或八叉树的一个问题是，单个叶节点可能包含在树的多个节点中。如果物体穿越了四叉树边界，它需要在两个节点中重新表达。球形树不具有这个性质。没有叶节点可以包含在多个球形树中。若球形树为真，则所有子节点也为真。如果球形树在视锥之外，则其所有的子节点也在视锥之外。对于范围测试和射线跟踪也是如此。

这使得球形树成为这类查询的理想数据结构。当通过球形树执行可视剔除时，每个节点都跟踪它在前一帧中的状态。仅当节点经历状态改变时，才发生回调（Callback），这能有效地保持仅存在视锥中物体的列表。

### 4.3.3 演示应用程序

Windows 应用程序 SphereTest.exe 中演示了这个算法。SphereTest.exe 将创建一棵包含 1000 个运动球体的球形树。尽管是以 2D 演示此应用，但球形树是一个完全 3D 数据结构。显示庞大球形树的帧频很低，因为在 Windows 下渲染所有的数据相当慢。球形树应用演示建立一个球形树，并执行各种不同的高速查询，用以模拟一个游戏中的典型场景。

模拟创建的球体个数可以作为命令行参数传递。通过较少的球体，更容易显现构建的球形树。如果创建 5000~10 000 个项的球形树，当大部分的 CPU 时间用于调用 Windows 图形程序来渲染结果时，查询速度仍然很快。如果应用大量项，当构建和平衡初始树时，可能会出现暂停，之后，运行速度也相当快。



这个模拟示例为实际游戏中可能看到的典型场景模拟。在这个模拟中，25% 的物体处于停止状态，其他 75% 的物体总是试图向 16 个吸引点之一发展。在实际游戏中，物体并不总是均匀分布在地址空间中。模拟示例演示使用包含在文件 Circle.cpp and Circle.h 中的 SpherePackFactory 类（请见附带光盘）。

图 4.3.1 显示了 SpherePack 系统的类图示。

要在模拟中使用 SpherePack 系统，只要用球体最大数、根节点大小、叶节点大小以及每个 SuperSphere 周围的增益量来实例化 SpherePackFactory 类即可。余量因子作为坐标空间中的溢出范围，避免子体与它们父体的 SuperSphere 分离太频繁。对模拟中的每个物体，调用 AddSphere() 方法来创建 SpherePack 实例。无论何时，若物体更改位置，就调用 NewPos() 方法。如果物体既改变位置，也改变半径，则调用 NewPosRadius() 方法。当销毁模拟中的物体时，在 SpherePackFactory 类上调用 Remove() 方法。



## SpherePack 系统

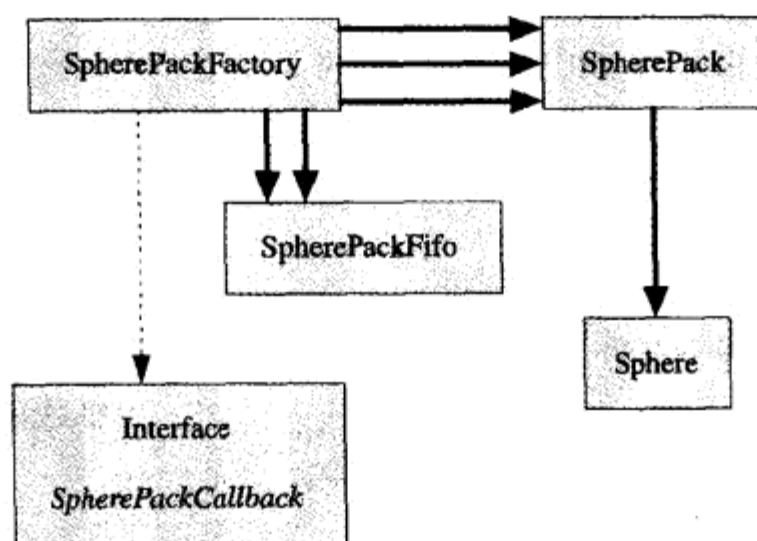


图 4.3.1 SpherePack 系统的类图示

**SpherePackFactory** 类保持完整的 **SphereTree** 的层次，而且包含两个集成的 **FIFO** 队列。当执行查询时，用一个叫 **SpherePackCallback** 的接口从 **SphereTree** 中提取信息。叶节点 **SpherePack** 继承了 **Sphere** 类的属性。



## 4.4 压缩的轴向包围盒树

Miguel Gomez, Lithtech Inc.

miguel@lithtech.com

**轴**向包围盒 (axis-aligned bounding box, AABB) 树结构被证明对于加速查询几何体集合的交集很有用。这个数据结构容易实现, 生成的结构很好地数字化了, 与所有的二叉树一样, 它们搜索时间为  $O(\log n)$  [Sedgewick90]。本文讲解了若干种技术, 可用于降低 AABB 树的整体内存消耗, 达到每个三角形 11 个字节。

### 4.4.1 概览层次排序方法

这部分内容包括四叉树、k-d 树、BSP 树、包围体积 (bounding volume) 树以及轴向包围盒, 它们是排序三维三角形集合的数据结构。

#### 1. 四叉树、k-d 树、BSP 树

一些排序三维三角形集合的最流行数据结构是八叉树、k-d 树和 BSP 树。虽然八叉树可能最容易实现, 但更复杂的 k-d 树和二分空间划分 (BSP) 树更适宜于特定的几何体集合。所有这些结构都使用平面将空间分成凸的区域, 这将导致一些严重的负面效应。首先, 跨在分割平面上的任何三角形要么被分成两块, 这就增加了数据集合, 并可能产生有问题的“长条”; 要么包含在多个节点中, 这将打破节点结构的上限。其次, 八叉树和 k-d 树不会适当地分割几何体, 因此, 它们可能需要一个任意阈值来中止递归构建算法。例如, 如果许多三角形共享相同的顶点, 则八叉树递归可能永远不能达到每节点一个三角形, 必须基于树深度来终止。这将导致不必要的深度树。适当地使用单个三角形的面进行分割的 BSP 树没有这个问题, 但它们仍然带来分割三角形的一些问题。

#### 2. 包围体积树

包围体积树解决问题的方式不同。包围体积树不是划分空间, 而是递归性地将三角形集合分解成两个子集, 并查找一个包含每个子集的包围体积 (bounding volume)。这种途径避免了必须切割三角形或在多个节点中包含它们, 而且当子集只有一个剩余三角形时终止树的构建过程。有关包围体积树的更多信息, 请参见参考文献[vandenBergen99]。

### 3. 轴向包容盒

轴向包容盒有一个定位（其中心）及范围。正如它的名字所表明的，AABB 的边平行于  $x$ 、 $y$  和  $z$  轴。图 4.4.1 显示一个 AABB 是如何包容单个三角形的。

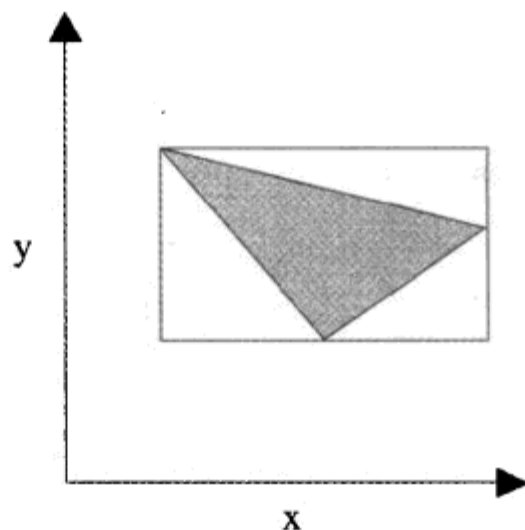


图 4.4.1 AABB 的边平行于坐标轴

### 4.4.2 AABB 树

AABB 中的边界框可以重叠，也经常重叠。图 4.4.2 显示一棵两个相连三角形（L 和 R）集合的小树。请注意，左 AABB 与右 AABB 重叠。

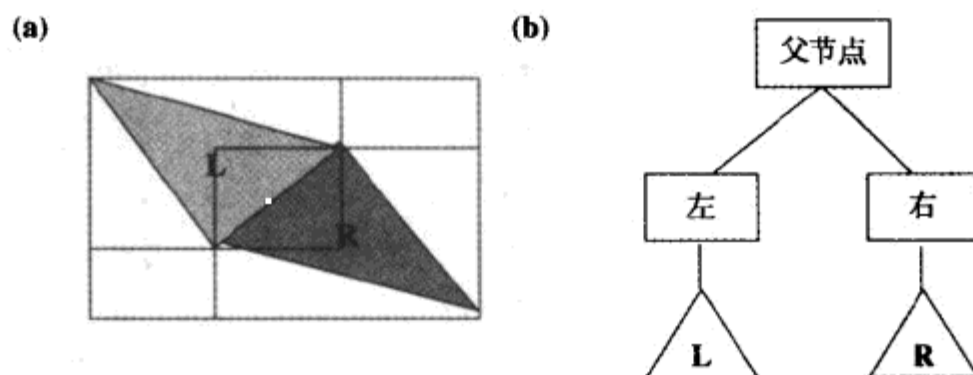


图 4.4.2 (a) 两个相连三角形的 AABB 总是相交。(b) 这个特定几何体的 AABB 树结构

### 4.4.3 构建 AABB 树

通过连续不断将三角形集合分解成两个子集，直到每个子集只包含一个三角形，此时就构建了一棵 AABB 树。构建 AABB 树的基本算法是：

- 查找包含整个集合的 AABB。
- 将集合分解成两个子集。沿着 AABB 的主轴（最长轴），按三角形的质心划分可以完成此任务。如果一个子集为空集，则任意创建两个大小近似相等的子集。
- 如果子集只包含一个三角形，则创建一个叶节点。
- 否则，对每个子集重复以上过程。

4.4.4 压缩 AABB 树

在大部分应用中，指定盒的范围，4 个字节的浮点数具有足够的精度了（7 位有效数字），每个节点需要  $6 \times 4 = 24$  个字节。任何二叉树刚好需要  $2n-1$  个节点来排序  $n$  个元素，在构建树前可以分配节点的数组。如果数据集庞大，则必须用 4 字节的无符号整数指定子节点的索引。为了避免对每个节点的每个三角形索引都必须存储另一个 4 字节的整数，可以利用第 0 个节点（即根）永远不会被另外的节点引用这一个事实。所以，如果其中一个子节点的一个索引为 0，可以将该节点看作一个叶节点，那么另外一个整数肯定指向一个三角形。根据该技术，每个节点共 32 个字节，或  $n$  变大时每个三角形约 64 个字节。如果想将三角形数量限制到  $2^{15}$ （节点数为  $2^{16}-1$ ），则对子节点索引可以使用无符号 2 字节整数。这将开销减少到每节点 28 个字节，或者每三角形 56 个字节。因为大部分应用可以将三角形集合分解成大小为  $2^{15}$  的块，所以可以将每三角形 56 个字节作为参考值。

AABB 有两个重要性质可用于最小化存储空间。首先，子 AABB 节点被它的父亲完全包含。根据此性质，可以将子节点范围（extent）值作为相对父节点范围偏移量来储存为一个无符号 8 位整数。其次，子节点最多有 6 个范围值不同于父节点。这意味着，完全描述两个子节点只需要保存 6 个值。

4.4.5 近似范围

为了把无符号字节用作节点的相对范围值，必须在构建树时，跟踪父节点的浮点数范围值。相对于父节点 AABB 换算出子节点范围，然后取最近整数值，得到子节点 AABB 范围的保守估计值（如图 4.4.3 所示）。

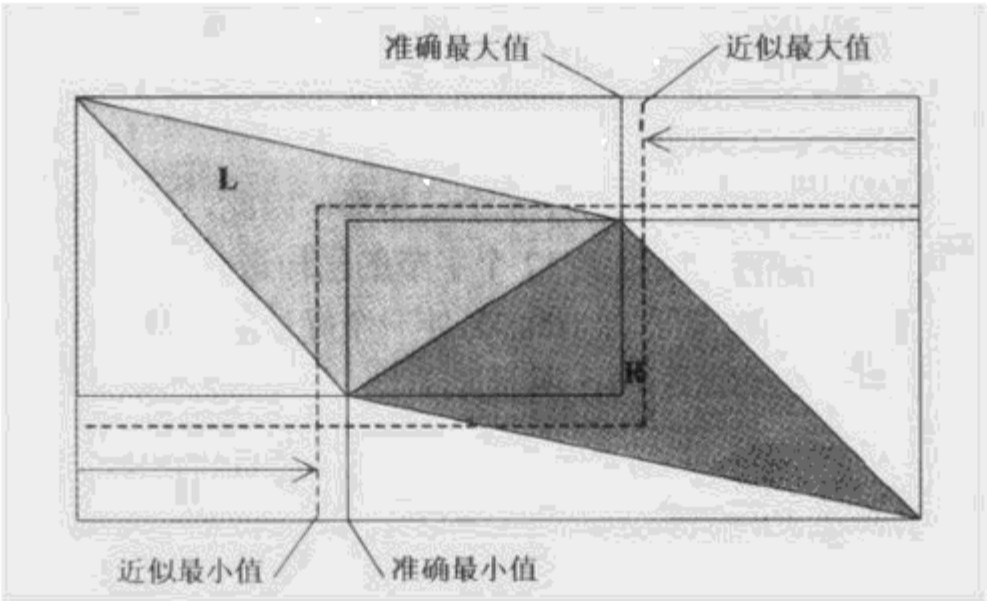


图 4.4.3 当通过无符号 8 位整数求近似范围值时，从左开始度量得到最小值，从右开始度量得到最大值

相对于父节点的所有范围值误差不会超过  $1/255 \approx 4 \times 10^{-3}$ 。按此规律，如果父节点的范围为每一边为 1m，则子节点的范围误差约为 4 mm。请记住，当节点变小时，以上数字将按比例变化。计算相对偏移值的公式是：

$$i_{min} = trunc \left( 255 \frac{\min_{child} - \min_{parent}}{\max_{parent} - \min_{parent}} \right)$$
$$i_{min} = trunc \left( 255 \frac{\max_{parent} - \max_{child}}{\max_{parent} - \min_{parent}} \right)$$
$$i_{max}, i_{min} \in [0, 255]$$

使用 1 字节取代 4 字节来保存范围值，内存消耗可以立即下降到每节点 10 个字节，或者每三角形 20 个字节。

4.4.6 利用冗余

为了对两个子节点只保存 6 个范围值，必须采取大胆措施，将左右子节点信息保存在一个节点结构中。在这个结构中，用一个位来表征哪一个范围属于哪一个子节点。真（true）位表示这个惟一的范围值属于左子节点（因此，父节点的范围值属于右节点）；假（false）位表示相反情况。这只需要 6 个位，因此其他 2 个位用于表示哪个节点为叶节点。当两个子节点的信息保存在单个节点中时，该节点只占用 11 个字节（如图 4.4.4 所示）。

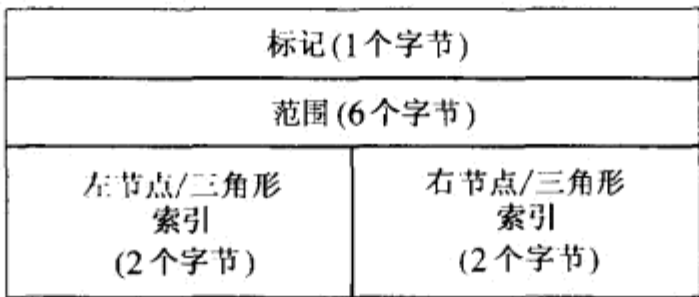


图 4.4.4 压缩 AABB 节点结构中，每节点只消耗 11 个字节

由于只需要  $n$  个节点（包括根）来排序  $n$  个三角形，因此，这棵树现在只需要每三角形 11 个字节的存储空间（一些架构可能需要 12 个字节的结构进行排列），这比单独使用相关值几乎少 50%，约为参考值的存储空间 的 1/5。另外一个好处是，可以跟踪双倍数量的三角形（即  $2^{16}-1$  个，而不是  $2^{15}$  个）（如图 4.4.5 所示）。

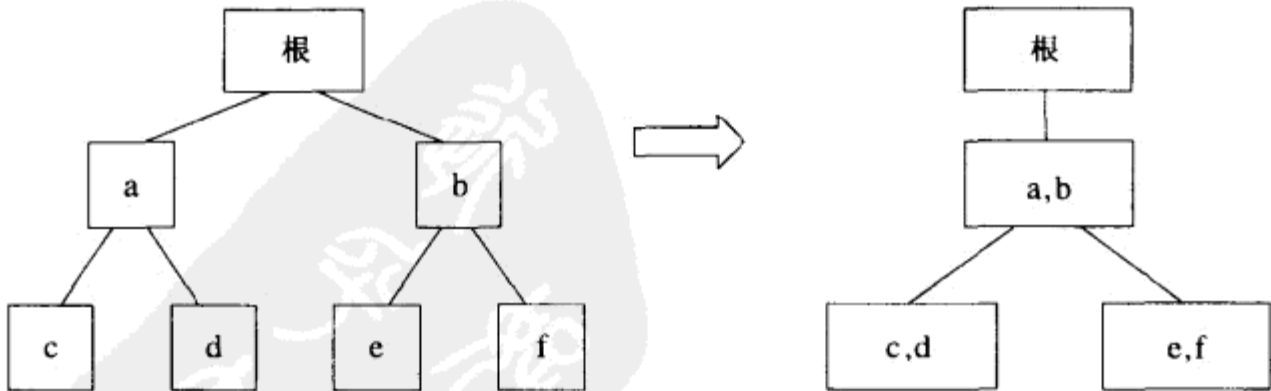


图 4.4.5 对于 4 个三角形的集合，未压缩树需要 7 个节点，而右边的树只需要 4 个节点

实际上,根节点是一个指定整个三角形集合浮点范围的单独数据类型,而且不是实际节点数组的一部分,它包含  $n-1$  个节点。

#### 4.4.7 运行时效

---

将 8 位值转回到浮点值的运行时开销很少。通过测试一条线段与 16K 随机分布的三角形的可能交集的试验表明,压缩结构约比完全浮点表达方式慢 10%。这将大大缩短数据组装入缓存的运行时间,几乎可以弥补把整数转换成浮点值的开销。

#### 4.4.8 将来的工作

---

相对于其他结构,AABB 树的压缩性具有明显的优势。不过,也可以应用相同的技术来压缩其他包围体积树结构,如 OBB 树。子 AABB 节点完全包含在父节点的体积中,根据此性质,可以缩放和共享范围值。此性质对 OBB 却不成立[Gottschalk96]。由于该性质不成立,因此不可能共享冗余的范围值。但是,从根节点扩展尺寸,繁殖成树,OBB 树仍然可以保存缩放的 8 位范围值以及 8 位方向值。

#### 4.4.9 参考文献

---

[Gottschalk96] Gottschalk, S., Lin, M. C., and Manocha, D., "OBBTree: A Hierarchical Structure for Rapid Interference Detection," SIGGRAPH Proceedings, 1996, pp. 171~180.

[Sedgewick90] Sedgewick, Robert, *Algorithms in C*, Addison-Wesley, 1990.

[vandenBergen99] van den Bergen, G., *Collision Detection in Interactive 3D Computer Animation*, Ph.D. Thesis, Eindhoven University of Technology, 1999.



## 4.5 直接访问四叉树查找

Matt Pritchard, Ensemble Studios

mpritchard@ensemblestudios.com

四叉树是从空间上管理大量二维数据的众所周知的通用构造。在游戏中，它们是控制游戏世界中物体位置的首选，它提供了优秀的搜索性能和灵活的实现可能性。在游戏数据为静态的场合，四叉树可以预计算，而且从性能的角度看，它们相当有效。

但对于动态性极强的场合，如在实时策略（RTS）游戏中，大量的物体不断地移动，在游戏世界（对应四叉树）中物体频繁的进入与退出，此时情况又如何呢？在实现以及后续剖析此场景的过程中，我发现四叉树更新程序非常惊人地消耗 CPU 的使用时间。这种不断更新数据集的行为使我决定去探寻这些时间到底花哪儿了？同时探寻更有效的途径来完成相同的工作。探寻的结果就是本文所讲解的内容。

该技术是对四叉树访问的一种通用优化技术，它针对传统实现从三个方面提高性能：（1）消除节点遍历以及它导致的不必要的缓存失败；（2）相对传统最佳情况，总体上要求较少的指令和 CPU 周期；（3）不管目标节点处于树中的深度如何，用最少的时间和最少的内存访问来完成。因为它的通用性，在很多特殊场合，可以发现许多方法改善这个技术。

### 4.5.1 性能剖析

检查许多原始四叉树遍历程序（与程序清单 4.5.1 类似）可知，经过编译器优化后，几乎不存在臃肿的代码。搜索程序只是几行检查每个子节点包围盒的代码，而且在必要时递归调用自身。即使在每次递归到四叉树底层的最坏情况下，生成调用次数的 CPU 周期计数估计值，它的总数甚至达不到剖析所指示的时间。

情况表明，性能拖累是现代计算机架构产生的结果。代码通过指针不断接触数据，然后移到内存中不相邻的其他数据。结果是产生相当高的缓存失败率。如果所有的四叉树更新按顺序排列，大部分朝向四叉树根节点的上部节点最终位于 CPU 缓存中。当更新没有排列在一起时，前面缓存的许多四叉树节点数据将在两次调用中刷新。

程序清单 4.5.1 中的代码包含了几乎整个比较函数，编译后将导致大量的分支、跳转和调用指令，而且对每一次循环存在大量可能的路径，这进一步降低了其性能。在现代 CPU 上，每次不希望的分支和跳转都会延迟指



令管道线。由于数据的性质，它的发生可能比想像的更频繁，因为与大部分循环不一样，前一个迭代中的分支选择不适用下一次选择。

### 4.5.2 消除中间阻碍

基本上说，程序清单 4.5.1 中紧凑的 C++ 代码花费在等待数据和指令的时间比实际执行的时间更多。进一步剖析后得知，大量被扫描的数据，实际并不需要，应该不在考虑之列。因此，问题就变成了在不实际检查四叉树节点本身的情况下，如何得到最终四叉树节点？

也就是问，或者更确切的是回答消除四叉树访问程序中的大部分性能障碍的问题。当然，实现计数（即本文所讲解的）也是回答该问题的一种有趣且非常有效的途径，但它决不是惟一可能的答案。

### 4.5.3 条件和要求

这个技术有两个限制。首先，四叉树必须规则，也就是说，每个节点正好在每个轴的中间分裂，让树中特定深度的所有节点表示相同大小的区域。幸运的是，大部分四叉树就是按这种方式实现的。其次，程序必须事先选择四叉树最大深度（子部分的数量）。同样，这也不是一种特例。考虑这些限制后可以认为，四叉树表示的面积必须为正方形，但情况并非完全如此，原因将在后面讲解。

惟一的要求是，程序需要让数组有序访问树中每一层的节点。这意味着，所有的四叉树节点指针放在空间位置排序线性数组中，或者在四叉树节点元素自身大小固定时，将它们保存在一个数组中。另一种方式是，给出特定层上的  $x$  和  $y$  节点坐标，四叉树数据的地址可以在不必遍历其他节点的情况下获得。

### 4.5.4 判断树层

直接查找目标四叉树节点的第一步是判断目标所在的树层。一旦知道了此值，搜索目标的坐标将提供直接查找目标节点的行与列。在余下的技术文章中，当我们谈到物体时，通常指的是它们的包围体积（bounding volume）。

如果它的体积与其中一根轴的节点中线相交，则认为物体位于四叉树的根节点。如果没有触及中线，则完全包含在一个子节点中（如图 4.5.1 所示）。就认为它已经通过了根体积包含测试。否则就找到物体与中线相交处的节点，重复以上过程，直到达到四叉树的底部。面积检查永远不必担心外部边，因为它们是前一层的中线和其他边。

现在来看代表四叉树面积的一根轴，假定其长度恰好是 2 的幂，比如为 256.0，四叉树有 9 层，其中最低层叶节点为  $1.0 \times 1.0$  的面积。根节点的中线为 128.0，下一层节点的中线为 64.0 和 192.0。3 层的中线为 32.0, 96.0, 160.0 和 224.0。看出它们之间的规律了吗？每一层的过渡点为 2 的不同幂，它们直接对应于树的层深度。

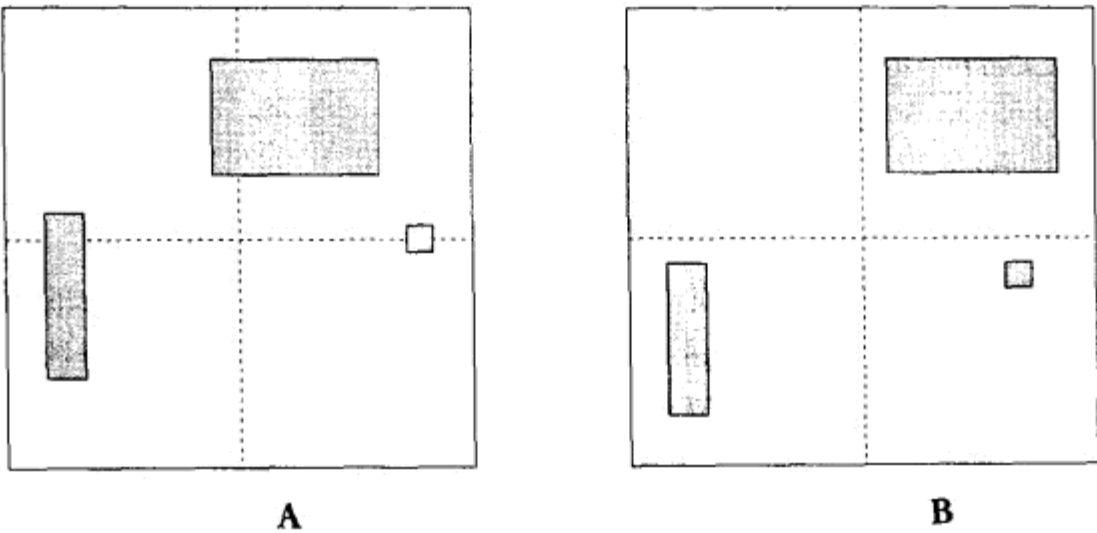


图 4.5.1 A)与节点中线相交的物体不能保存在子节点中；B)不与中线相交的物体保存在子节点中

现在取搜索物体包围面积的相同轴，得到一条线段 $(x_1, x_2)$ 。我们希望知道它跨越的最大 2 的幂，因为这表示物体可以完全包含于四叉树最低层。“跨越”2 的幂只是意味着从低于 2 的幂值到大于（或等于）该值， $x_1$  到  $x_2$  的过渡范围。因为存储要求对象完全包含在四叉树节点内，所以保存物体的四叉树层次实际上为更高一层，要高于物体跨越的节点面积由特定 2 的幂表示的那一层。

因为值的跨越可以表示为 0 到 1 的二元过渡，所以程序求范围 $(x_1, x_2)$ 的整数部分并将它们异或（XOR）在一起，可以快速判断跨越发生的地方。结果中的每个位表示与 2 的幂对应的点相交的范围，如果在该位置有一个过渡，则返回 1。结果中首个 1 位的位置越低，跨越在四叉树中可以保存得越低（越深）。如果出现 0 结果（即结果中无 1 位），则范围可以保存在四叉树中非常低的层。因此，最高“1”位的位置表示从四叉树底层开始，可以在多高的层次上首次正确安放范围 $(x_1, x_2)$ 。

下面是使用 9 层四叉树长度轴为 256.0 的例子，演示这个 2 的幂的性质。图 4.5.2 显示每个例子的 2D 表达方式，使用的是  $x$  轴。

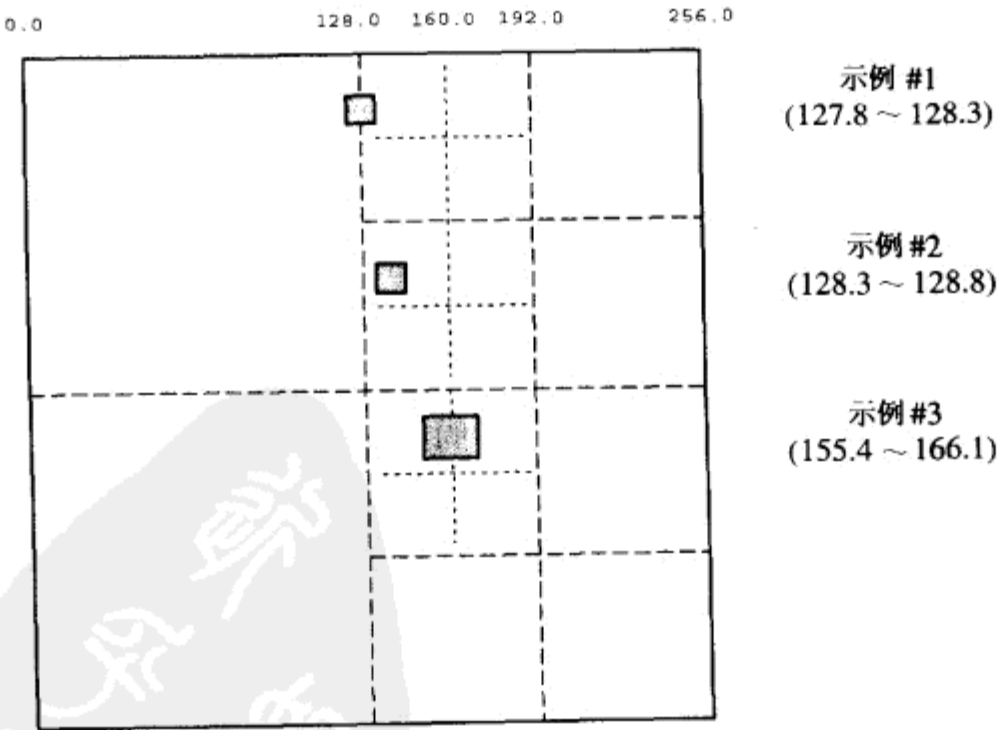


图 4.5.2 示例的图示。为了演示的目的，没有度量示例框

## 示例 1

物体包围(127.8, 128.3)。这物体位于根节点的中间, 仅适配于顶层。 $(x_1, x_2)$ 的整数部分为(127, 128)。

```
X1=      127 = 01111111
X2=      128 = 10000000
-----
XOR result 255 = 11111111
```

结果中的第一个1为位7( $2^7$ ), 即第8个位置。 $9-8$ (树层数减去位置)得到1, 即四叉树的第一层(根层)。

## 示例 2

物体包围(128.3, 128.8)。示例1中的物体被推离中线, 下落到叶节点(9层)。 $(x_1, x_2)$ 的整数值为(128, 128)。

```
X1=      128 = 10000000
X2=      128 = 10000000
-----
XOR result  0 = 00000000
```

结果中没有1, 因此, 得到位置0。四叉树的 $9-0=9$ 层。

## 示例 3

物体包围(155.4, 166.1)。这是一个非常大的物体, 中心位置约在范围的60%处。因为其大小相当于最小节点, 所以它至少是底层节点(9层)上的4层。 $(x_1, x_2)$ 的整数值为(155, 166)。

```
X1=      155 = 10011011
X2=      166 = 10100110
-----
XOR result  61 = 00111101
```

第一个1在位位置6。四叉树的 $9-6=3$ 层。

## 4.5.5 位置映射

现在, 我们知道如何判断一维树的层了, 将这一过程推广到二维, 只需要对每根轴重复这一过程, 并选择四叉树中最高层的结果。

虽然前一个示例能发挥作用(因为它们假定存储范围为2的幂), 但在大部分情况下, 四叉树表示的面积不是偶数2的幂, 而且最小节点也不是正好 $1.0 \times 1.0$ 的大小。所幸的是, 将坐标值从游戏或应用的度量系统按比例转换为四叉树的尺寸就可以解决。在初始化时计算每根轴的比例因子, 并在转换到整数前将坐标乘以因子就可以完成。对于四叉树比例, 1.0表示每根轴上叶节点的大小。(请注意每根轴可以有不同的比例因子, 允许将矩形面积映射为正方形节点)。因而,  $n$ 层四叉树的比例坐标中, 每根轴的范围在 $0 \sim 2^{n-1}$ 。

缩放和转换坐标到整数后，异或每根轴的起始和终止位置。判断哪个位置包含最高 1 位的方式有多种。最简单而且移植性最好的方式是循环移位计数位数，但对于具有 6 或更少层的四叉树，使用查找表可能更适合。更理想的是特定平台解决方案，如奔腾（Pentium）芯片上的 BSR 指令以及 PowerPC 上的 cntlzw 指令，它们位扫描某个值，去掉了该过程中的循环和分支。程序清单 4.5.2 是一个代码示例，它使用上述方法和一个 while 循环确定四叉树的层次。

#### 4.5.6 判断位置

---

一旦判断了四叉树中的层，剩下的一步是得到缩放过的坐标，并提取目标节点的行列位置。如果目标节点在底层，则不需要缩放。对于之上的每一层，目标节点的坐标必须乘以缩放因子 2，以反映该层上的较少节点数。也可以右移整数坐标值完成缩放。右移数为从底层开始到目标层的数目，这样得到行和列位置，将其插入寻找四叉树该层上节点的数组中。程序清单 4.5.3 演示了一个综合上述功能的示例函数。

#### 4.5.7 遍历四叉树

---

定位搜索节点后，如果仍然需要遍历四叉树，所要做的是保存树层和数组行与列位置。向左或向右移动位置值，程序可以向上移或向下移访问树层。同时，增加或减小数组位置，它可以移到当前树层邻层。

#### 4.5.8 优化四叉树

---

有两个技巧可以充分发挥四叉树的作用，这些技巧适用于任何四叉树访问实现代码。

第一个技巧是确保项放置在树中尽量低的位置。不正确的边界条件可能导致项放在比适当位置高得多的位置。在使用一棵 256×256 树的例子中，跨越位置 128.0 的物体将放在根节点。不过，对于那些只触及但未实际跨越的物体，如大小为 1.0、位于 127.0 到 128.0 的“分片”，情况又如何呢？在这种情况下，将 “>= ” 或 “ <= ” 转换为 “> ” 或 “< ”，或者少量移动物体坐标可以区分根节点位置与叶节点位置。

第二个技巧是，确保程序选出四叉树的正确层数。如果挑选出的太少，搜索项时，节点将被重载，此时全部思路是首先减少搜索集的大小。太多的树层数同样不好。许多节点可能为空，浪费了空间，而且在不同的运算中将出现太多的节点遍历。最好在四叉树实现的非发布版本中添加性能计数器及统计量，同时使用实是数据分析结果，因为最优设置取决于碰到的实际数据。

最后要说的是，这一技术应该也适用于八叉树。也许有人对八叉树尝试了这一技术，请让我知道它的工作性能。

## 程序清单 4.5.1 四叉树搜索的简单实现

```
// Two C++ classes, one to represent the QuadTree and one to
// Represent each individual node.
// all variables are class members unless defined

QuadNode* QuadTree::GetNodeContaining(const rect &ObjBounds)
{
    if (RootNode->Contains(ObjBounds)
        return(RootNode->GetNodeContaining(ObjBounds);
    else
        return(NULL);
}

QuadNode* QuadNode::GetNodeContaining(const rect &ObjBounds)
{
    if (!isLeafNode)
    {
        if (ULNode->Contains(ObjBounds)
            return(ULNode->GetNodeContaining(ObjBound));

        if (URNode->Contains(ObjBounds)
            return(URNode->GetNodeContaining(ObjBound));

        if (LLNode->Contains(ObjBounds)
            return(LLNode->GetNodeContaining(ObjBound));

        if (LRNode->Contains(ObjBounds)
            return(LRNode->GetNodeContaining(ObjBound));
    }
    return(this);
}

bool QuadNode::Contains(const rect &ObjBounds)
{
    return(ObjBounds.top >= y1 && ObjBounds.left >= x1 &&
        ObjBounds.bottom <= y2 && ObjBounds.right <= x2);
}
```

## 程序清单 4.5.2 判断目标四叉树节点的层的代码

```
int QuadTree::GetNodeLevelContaining(const rect &ObjBounds)
{
    int xResult = ((int) (ObjBounds.left * QuadXScale)) ^
        ((int) (ObjBounds.right * QuadXScale));

    int yResult = ((int) (ObjBounds.top * QuadYScale)) ^
        ((int) (ObjBounds.bottom * QuadYScale));
```

```

int NodeLevel = NumberOfTreeLevels;

while (xResult + yResult != 0 )    //Count highest bit position
{
    xResult >>= 1;
    yResult >>= 1;
    NodeLevel--;
}

return (NodeLevel);
}

```

#### 程序清单 4.5.3 判断目标四叉树节点的层的代码

```

QuadNode* QuadTree::GetNodeContaining(const rect &ObjBounds)
{
    int x1 = (int) (ObjBounds.left * QuadXScale);
    int y1 = (int) (ObjBounds.top * QuadYScale);

    int xResult = x1 ^ ((int) (ObjBounds.right * QuadXScale));
    int yResult = y1 ^ ((int) (ObjBounds.bottom * QuadYScale));

    int NodeLevel = NumberOfTreeLevels;
    int shiftCount = 0;

    while (xResult + yResult != 0 )    //Count highest bit position
    {
        xResult >>= 1;
        yResult >>= 1;
        NodeLevel--;
        ShiftCount++;
    }

    // Now lookup the node pointer in a 2D array stored linearly

    x1 >>= shiftCount;                // Scale coordinates for
    y1 >>= shiftCount;                // quadtree level

    QuadNode** nodes = NodeLevelPointerArray[NodeLevel];

    return (nodes[y1<<(NodeLevel-1)+x1]);
}

```





## 4.6 近似鱼缸折射

Alex Vlachos, ATI Research

alex@vlachos.com

这篇技术文章简短讲解一种方法,用于近似在鱼缸边看到的反射效果。重点是解释如何构建转换矩阵,应用于鱼缸内的几何体,来模拟折射效果。

### 4.6.1 鱼缸观察现象

您有过在海鲜餐馆长时间等待,久久盯着大厅旁边鱼缸的经历吗?如果有,您大概注意到了鱼缸内部形状改变的独特方式。通过观察,您将发现,在鱼缸旁边从一边到另一边走动,鱼缸的后壁仿佛在左右方向发生了错切。同时,当站在鱼缸左边远处或右边远处时,鱼缸似乎扁平成一幅 2D 图像。进一步观察表明(这也是常识),观察者头部的高度也相似地影响鱼缸的内部形状。可以通过结合错切和缩放矩阵的方法来模拟以上现象。

虽然使用 Snell 定律,可以用一些在数学上更准确的方法来计算折射光,但本文的重点只是基于经验观察获得近似的视觉效果。务必注意,本文讨论的是鱼缸,而不是鲨鱼缸,两者之间的区别在于相对于观察者,鱼缸的大小不同。这个算法假定一般情况下,观察者离玻璃前端不是特别近。“特别近”显然是一种相对而言的措词。如果鱼缸尺寸约为 3 英尺×2 英尺,观察者在附近走动,这正是本文的讨论情况。如果同一位观察者在 30 英尺×10 英尺的鲨鱼缸前面走动,可以说,这个人大部分时间位于鱼缸的正前方,离玻璃相当近。虽然这时算法也不会失败,相对于玻璃运动而不是站立在正前方时,才能观察到最佳视觉效果。借鉴到游戏环境中,也有相似的假定,这一点对于实时图形没有什么不同。

#### 1. 预计算

已知一个鱼缸,它的前端尺寸为 width (宽) 和 height (高),将鱼缸前表面中心放原点的转换矩阵可以预计算(例如,鱼缸后部沿正 z 轴扩展,将正 y 轴作为向上的矢量)。前表面的中心放在原点的原因是便于缩放和错切转换。一旦产生了转换矩阵,摄像机的位置可借助于乘以简单矩阵转换到鱼缸空间。这将允许所有的计算在此公共空间(鱼缸空间)内发生。需要注意的是,这个方法仅基于摄像机相对于鱼缸的位置,摄像机的观察方向不影响视觉效果。

## 2. 缩放因子

对每帧都需要计算的鱼缸缩放因子，可以形象表达为将鱼缸从远处表面朝近处表面（即通过该面观察）移动。当观察者站在玻璃正前方，则不发生缩放。当观察者站在鱼缸的侧面透过玻璃表面观察时，则大比例缩放鱼缸。从左到右移动的效果如图 4.6.1 所示。

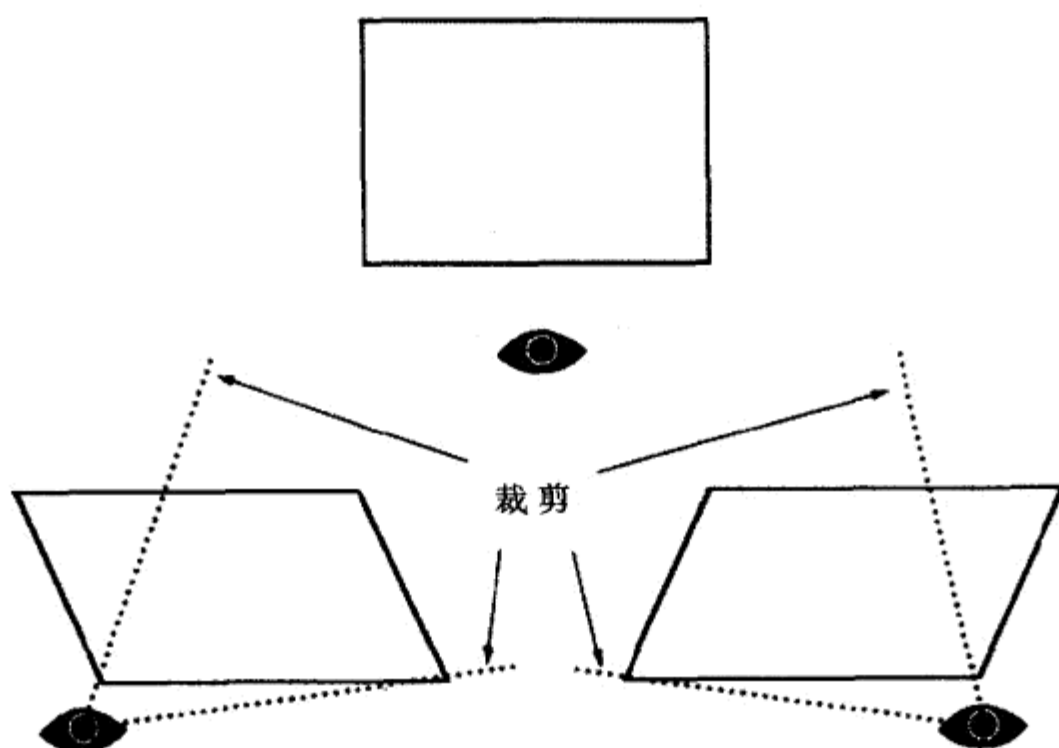


图 4.6.1 基于观察者位置的 3 个错切和缩放示例。用户裁剪面基于观察者位置和鱼缸玻璃的边创建

已知鱼缸空间中观察者的位置以及玻璃的中心，计算鱼缸表面的角度时应考虑鱼缸玻璃的宽高比值。如果不考虑玻璃尺寸，透过不同玻璃形状的结果就不一致。

使用如下方法计算缩放因子：

- (1) 将观察者的位置转换到鱼缸空间。这把将鱼缸前面玻璃的中心放置在原点。
- (2) 计算从玻璃中心（原点）到摄像机的矢量，但不归一化（normalize）。
- (3) 将矢量的  $x$  和  $y$  坐标分别除以玻璃的宽和高。
- (4) 将矢量的  $z$  坐标除以鱼缸的深度（由近到远的距离）。
- (5) 单位化矢量。
- (6) 计算矢量与鱼缸空间  $(0, 0, -1)$  中玻璃的面法线的点积（dot product）。钳制 0.0 到 1.0 的标量值作为缩放因子。

假设按前面的解释，将鱼缸转换到原点，缩放沿着  $z$  轴进行。

## 3. 错切因子

错切因子也必须对每一帧进行计算。需要生成两个单独错切值：一个针对  $x$  轴，另一个针对  $y$  轴。为了沿着  $x$  轴错切，创建一个由  $x$  和  $z$  坐标组成的 2D 矢量并归一化，坐标系即用来计算机缩放因子的坐标系。单位化 2D 矢量的  $x$  分量用于沿着  $x$  轴的错切。与之相似，创建由  $y$  和  $z$  坐标组成的 2D 矢量并归一化，并使用归一化的 2D 矢量的  $y$  分量用于沿着  $y$  轴的错切。试验表明，将  $x$  向的错切值限制为 0.75， $y$  向的错切值限制为 0.25，有助于减少不希望

望出现的效应。

#### 4. 组成矩阵

缩放因子和错切因子可以放入一个矩阵中：

$$\begin{bmatrix} 1 & 0 & shear_x & 0 \\ 0 & 1 & shear_y & 0 \\ 0 & 0 & scale & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

#### 5. 用户剖面 (Clip Plane) 和模板缓冲 (Stencil Buffer)

当应用错切和缩放转换几何体后，可能有些内容不需要更改。如果鱼缸是矩形，可以通过产生从观察者位置到玻璃边界的4个用户剪裁面来解决这个问题（如图4.6.1所示）。这将剪裁所有从鱼缸外部可以看到的错切几何体。如果玻璃为环形或其他的非矩形，则需要用模板缓冲替换用户裁剪面，来剔除不属于鱼缸玻璃的像素。

### 4.6.2 提高真实性

到目前为止，这个方法描述了如何让鱼缸里面的几何体看起来像在充满了水的鱼缸中一样。还可以进行更多的改进。第一个改进是应用刻蚀纹理 (caustic texture)。简单地说，刻蚀是出现在充满水的容器底部的一种光线图案。当光从鱼缸外部穿过顶部水表面时，光的折射将形成这些图案。既然使用刻蚀效果可以不错地模拟水的表面，也可以通过投影纹理使鱼缸中所有几何体上滚动两个刻蚀纹理，以产生刻蚀效果。

下一个明显可以改进的是让玻璃外部具有反射效果。一种得到玻璃上理想反射的简单途径是将反射的玻璃几何体浸染成一个可渲染纹理，然后进行纹理渲染，得到透明反射表面的幻觉。

为了透过鱼缸顶部的水进行观察，可以先使用一个可渲染纹理来渲染假如鱼缸中没有水可以看到的物体，然后应用环境反射凹凸贴图到那个可渲染纹理上以得到动画水效果（或其他像素着色器）。

显然，这个算法可以应用到鱼缸的每个单独表面。另外，如果定义玻璃每边的顶点在当前视锥 (view frustum) 的外部，则可以执行更高级别的裁剪。裁剪庞大多边形组的方法适用于玻璃的每个平面。

### 4.6.3 结论

可以很容易地制造出平面的模拟折射效果的假像。结合其他实时渲染技术，可以创建非常逼真的效果。

## 4.7 渲染打印分辨率的屏幕快照

Alex Vlachos , Evan Hart, ATI Research

alex@vlachos.com , ehart@ati.com

由于最新实时图像技术的发展, 游戏图像的质量得到了极大的提高。尽管在计算机显示器上 (一般分辨率为 72dpi), 这些图像的显示效果已经相当不错, 但这个分辨率对打印来说显得有点偏低。打印分辨率一般是显示器二维分辨率中较小者的 4 倍。因此, 在打印页面上, 将屏幕像素按“一对一”关系映射到打印机的点, 导致屏幕快照将只有邮票大小。解决这个问题的一种方案是, 在页面排版软件中扩大原始图像, 但这将导致斑驳的图像, 不能充分表现出原始图像内容。

这个问题的简单解决方案是, 从更高分辨率的帧缓冲获取屏幕快照。一般而言, 更大的缓冲区可以超出屏幕显示区的表面, 因而不受显示器所能支持分辨率的限制。但这仍有一些问题, 因为大部分图形加速器设计成仅能渲染显示器的最大分辨率, 使得程序员最高只能获得  $2048 \times 2048$  分辨率的屏幕快照。尽管这个分辨率看来已经足够了, 但尚不能全部覆盖 300dpi 分辨率的杂志页面大小。使用这种方法, 每种颜色缓冲和深度缓冲需要 32MB 的显存。

本文提出的解决方案是, 将该任务分解成渲染一些更小的子图像。这些子图像可以被粘接到一起, 组成任意分辨率的屏幕快照。本文技术的重点是无缝连接子图像所需要的投影矩阵 (projection matrix) 技巧。

### 4.7.1 基本算法

既然现有的图形硬件无法渲染打印要求的分辨率, 最终期待的图像必须按方格划分成一系列的子图像渲染, 然后拼贴在一起。在本文后面给出的示例代码中, 假定每个子图像的尺寸与帧缓冲尺寸相等。尽管这些子图像的分辨率必须被图形硬件支持, 但是将视口尺寸设置成小于帧缓冲的某个值, 这个方法就可以处理任意尺寸的图像。

图 4.7.1 显示如何用网格将图像分割成一系列子图像。6 个视锥面定义了一个投影矩阵: 近、远、左、右、顶、底视锥面。对于每一个子图像, 计算新的左、右、顶和底视锥面将定义惟一的投影矩阵来渲染子图像。

获得屏幕快照分三步:

- (1) 计算视锥截面;
- (2) 决定中间截面;

(3) 使用由中间截面组成的每个子图像的视锥渲染子图像。

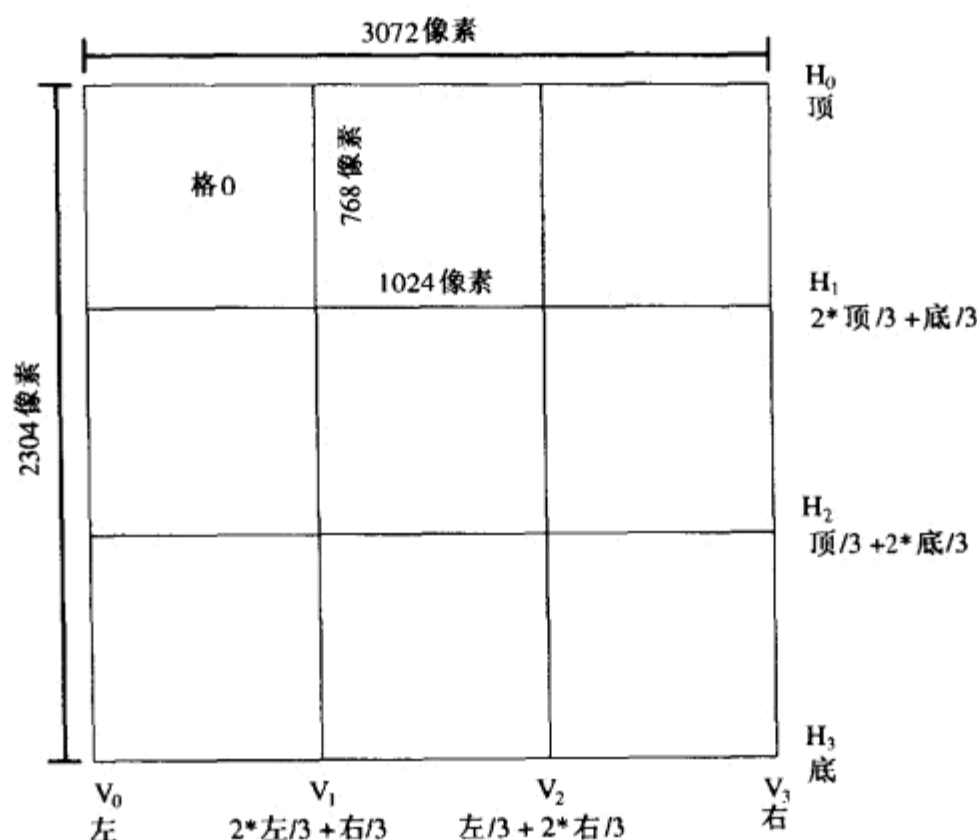


图 4.7.1 图像分割

首先，计算投影矩阵的视锥截面获得 6 个标量。近和远标量就是从眼睛到近、远平面的距离，同时左、右、上、下标量值由该平面与近平面相交的点定义。“实时渲染”一书的 3.5 节准确描述了这些值的涵义[Möller99]。

接下来，通过分别线性内插上、下值以及左、右值，来计算中间面值  $H_n$  和  $V_n$ 。

最终，每个子图像的视锥由左、右、上、下截面对应的边界  $H_n$  和  $V_n$  值组成。可以调用 `glFrustum()` 或者 `D3DXMatrixPerspectiveOffCenter()` 以获得投影矩阵。当每个子图像渲染完成，整个图像也就渲染完毕。

下面是人工构建  $3 \times 3$  栅格子图像的投 OpenGL 和 Direct3D 影矩阵的简单示例代码。

```
const float GPG_PI = 3.14159265f;
inline float GpgDegToRad(float D) { return ((D) * (GPG_PI/180.0f)); }
void GpgPerspective (double fovy, double aspect, double Near, double Far, int subrect)
{
    double fov2, left, right, bottom, top;

    fov2 = GpgDegToRad(fovy) * 0.5;

    top = Near/(cos(fov2)/sin(fov2));
    bottom = -top;
    right = top*aspect;
    left = -right;

    if (subrect == -1) //常规满屏 Regular full screen
        GpgFrustum (left, right, bottom, top, Near, Far);
}
```

```

else if (subrect == 0) //顶左 UL
    GpgFrustum(left, left/3.0, top/3.0, top, Near, Far);
else if (subrect == 1) //顶当前 UC
    GpgFrustum(left/3.0, right/3.0, top/3.0, top, Near, Far);
else if (subrect == 2) //顶右 UR
    GpgFrustum(right/3.0, right, top/3.0, top, Near, Far);
else if (subrect == 3) //中左 ML
    GpgFrustum(left, left/3.0, bottom/3.0, top/3.0, Near, Far);
else if (subrect == 4) //中当前 MC
    GpgFrustum(left/3.0, right/3.0, bottom/3.0, top/3.0, Near, Far);
else if (subrect == 5) //中右 MR
    GpgFrustum(right/3.0, right, bottom/3.0, top/3.0, Near, Far);
else if (subrect == 6) //底左 BL
    GpgFrustum(left, left/3.0, bottom, bottom/3.0, Near, Far);
else if (subrect == 7) //底当前 BC
    GpgFrustum(left/3.0, right/3.0, bottom, bottom/3.0, Near, Far);
else if (subrect == 8) //底右 BR
    GpgFrustum(right/3.0, right, bottom, bottom/3.0, Near, Far);
}

void GpgFrustum (double left, double right, double bottom, double top, double zNear,
double zFar)
{
    float matrix[16] = { 1.0f, 0.0f, 0.0f, 0.0f,
                        0.0f, 1.0f, 0.0f, 0.0f,
                        0.0f, 0.0f, 1.0f, 0.0f,
                        0.0f, 0.0f, 0.0f, 1.0f };

#ifdef GPG_OPENGL_API
    matrix[0] = (float)(2.0*zNear/(right-left));
    matrix[5] = (float)(2.0*zNear/(top-bottom));

    matrix[8] = (float)((right+left)/(right-left));
    matrix[9] = (float)((top+bottom)/(top-bottom));
    matrix[10] = (float)(-(zFar+zNear)/(zFar-zNear));
    matrix[11] = (float)(-1.0);

    matrix[14] = (float)(-(2.0*zFar*zNear)/(zFar-zNear));
#else //Direct3D
    matrix[0] = (float)(2.0*zNear/(right-left));
    matrix[5] = (float)(2.0*zNear/(top-bottom));

    matrix[8] = (float)((right+left)/(right-left));
    matrix[9] = (float)((top+bottom)/(top-bottom));
    matrix[10] = (float)(-(zFar)/(zFar-zNear));
    matrix[11] = (float)(-1.0);

    matrix[14] = (float32)(-(zFar*zNear)/(zFar-zNear));
#endif

    //现在设置该矩阵为当前投影矩阵
}

```



### 4.7.2 忠告及注意

---

对于任何渲染算法，程序员必须注意该技术如何与图形引擎执行的特殊技巧相互影响。程序员也应该明白，若场景中存在非常大的多边形，利用多图像形成屏幕快照可能会影响 Gouraud 着色 (Gouraud shading)。由于多投射，大多边形可能在所有子图像边界处被裁剪，这将导致光照不自然的效果。另外，程序员也应该小心处理需要渲染成纹理的特效。为了获得最佳效果，也需要提高这些图像的分辨率。当然，作为动画场景，程序员应该确保动画中各个子图像渲染的同步。

当我们直接引用 Direct3D 和 OpenGL 的图形 API，并且使用 PC 硬件时，即使在屏幕快照问题最不利的游戏控制台上，这项技术同样能很好地工作。若不使用这项技术在这些设备上，游戏杂志受到截屏问题的种种束缚，只能从 NTSC 信号中获取糟糕的图像。

这是一项欺骗技术吗？不，这种方法等价于在高分辨率下玩游戏。不会由于提供给游戏杂志高打印分辨率的屏幕快照图像而让别人误解您的工作。我们已经在 *Real-Time Rendering* 和 *Game Programming Gems* 这两本书的插图中使用了这种技术，获得了满意的效果。

### 4.7.3 结论

---

我们提供了使用图形硬件渲染打印分辨率屏幕快照的技术。由于不可能直接渲染符合打印要求的高分辨率图像，我们演示了如何将所期望得到的最终图像分割成可以单独渲染、按网格分布的子图像。这将使得游戏内容开发者能在打印介质上以最好的质量展示游戏。

### 4.7.4 参考文献

---

[Möller99] Möller, Tomas, and Eric Haines, *Real-Time Rendering*, AK Peters, 1999.



## 4.8 对任意表面应用贴花

Eric Lengyel, C4 Engine

lengyel@c4engine.com

许多游戏需要渲染墙上焦痕和地面脚印之类的特效，这些并非原始场景中的组成部分，而是在游戏进行中即时生成的。这些特效通常是通过创建一个新对象来实现的，我们称之为贴花，贴花与现有表面一致，使用某些类型的深度偏移技术（depth offset technique）渲染而成（示例请见参考文献[Lengyel00]）。对一个二维表面的内部进行贴花是容易的，但对目前游戏中用来表现弯曲物体和地形块等对象的复杂表面进行贴花时，困难就出现了。本文提供一种通用方法，可以对任意形状表面进行贴花，并沿着任意表面边界对贴花进行统一剪裁。

### 4.8.1 算法

从已知表面的一点  $P$  以及垂直该点处表面的单位法线方向  $N$  开始。点  $P$  代表贴花的中心，可能是子弹击中表面的点，或者游戏人物踩踏地面的点。同时必须选择单位切线方向  $T$ ，以决定贴花的方向。整个构造如图 4.8.1 所示。

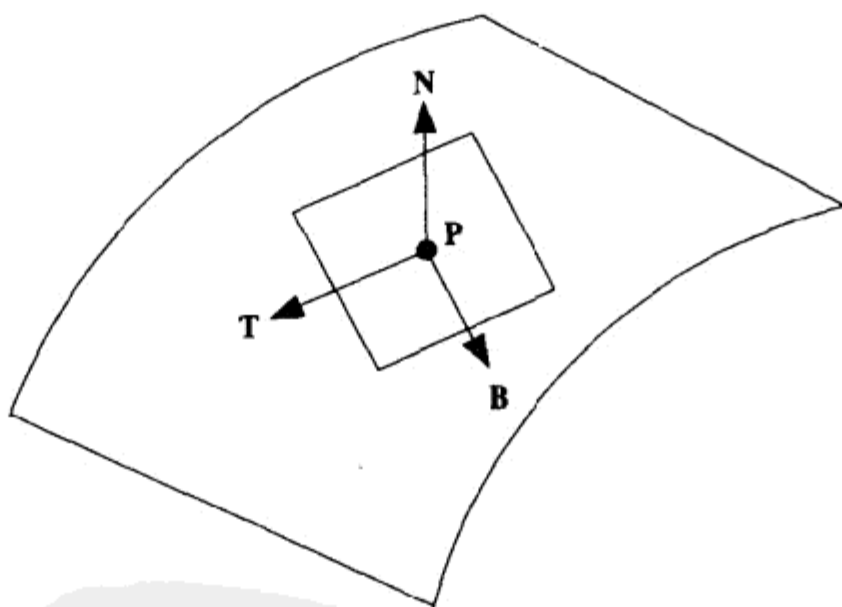


图 4.8.1 贴图框架

给定点  $P$ 、方向  $N$  和  $T$ ，就形成了在  $P$  点与几何体表面相切的导向平面。可以通过建立平行于法线方向  $N$  的 4 个边界平面，从导向平面中划定一个矩形，该矩形表示贴花区域。令  $w$  和  $h$  表示贴花的宽和高，则对应于

4 个边界平面的 4D 向量由公式 4.8.1 给出:

$$\begin{aligned} left &= \left( T, \frac{w}{2} - T \cdot P \right) \\ right &= \left( -T, \frac{w}{2} + T \cdot P \right) \\ bottom &= \left( B, \frac{h}{2} - B \cdot P \right) \\ top &= \left( -B, \frac{h}{2} + B \cdot P \right) \end{aligned} \quad (4.8.1)$$

其中  $B=N \times T$ 。我们将通过剪裁 4 个边界平面附近的表面, 为贴花对象生成一个三角形网格。同时也希望经过前后平面的剪裁, 避免穿透到可能位于同一边界平面内、但处于  $P$  点太前或太后的相同表面网格部分。对应于前后平面的 4D 向量由公式 4.8.2 给出:

$$\begin{aligned} front &= (-N, d + N \cdot P) \\ back &= (N, d - N \cdot P) \end{aligned} \quad (4.8.2)$$

其中,  $d$  是贴花中任意顶点与经过点  $P$  的切平面的最大距离。

算法过程如下所述。首先, 我们决定游戏场景中哪些表面可能受到贴花的影响, 以那些表面的包围体积与  $P$  点距离为定值作判据。对每一个可能受影响的表面, 分别检查表面网格中的每一个三角形。令  $M$  表示网格中三角形的单位法向。 $\epsilon$  为一确定正值。按  $M \cdot N < \epsilon$  的关系, 排除一些背离法线方向  $N$  的三角形, 保留由公式 4.8.1 和 4.8.2 给定平面的剪裁的三角形, 并将其存储在一个新的三角形网格中。

当一个三角形与任意平面重叠而需要剪裁时, 插值 (interpolate) 法向量以及顶点位置, 稍后可以在错切顶点上应用这些颜色值。这些颜色反映了每一个顶点法向与贴花法向之间的夹角。这将具有按照每一个三角形与贴花平面的相关方向平滑渐变贴花纹理的效果。应用公式 4.8.3 为每一个顶点指定一个  $\alpha$  值:

$$\alpha = \frac{\frac{N \cdot R}{\|R\|} - \epsilon}{1 - \epsilon} \quad (4.8.3)$$

其中  $R$  是顶点的法向 (可能由于插值而非标准化)。这样就将点积区间  $[\epsilon, 1]$  映射成  $\alpha$  值区间  $[0, 1]$ 。

通过度量经过  $P$  点的距离, 同时利用法线方向  $T$  和  $B$ , 将纹理映射 (texture mapping) 坐标应用到最终的三角形网格中。若  $Q$  表示贴花三角形网格一个顶点的位置, 则纹理坐标  $s$  和  $t$  由下式给出:

$$\begin{aligned} s &= \frac{T \cdot (Q - P)}{w} + \frac{1}{2} \\ t &= \frac{B \cdot (Q - P)}{h} + \frac{1}{2} \end{aligned} \quad (4.8.4)$$

### 4.8.2 三角形剪裁

曲面中每一个可能受到贴花影响的三角形都被作为凸多边形处理，依次裁剪到 6 个边界平面。剪裁具有  $n$  个顶点的凸多边形到某个平面将产生一个最多  $n+1$  个顶点的新凸多边形。这样的话，相对所有 6 个平面剪裁的多边形最多具有 9 个顶点。一旦剪裁完成，每个多边形被作为三角形扇区处理，并被加入到贴花的三角形网格中。

为了沿着任意平面剪裁凸多边形，首先将多边形的所有顶点分为两类：位于平面背面的属于一类，位于平面正面以及平面上的属于一类。如果多边形所有的顶点位于平面的背面，该多边形不予考虑。否则，我们查看多边形中两两相邻的顶点，寻找与裁剪平面相交的边。如图 4.8.2 所示，在相交的地方，在多边形中加入新顶点，删除裁剪平面背面的顶点。

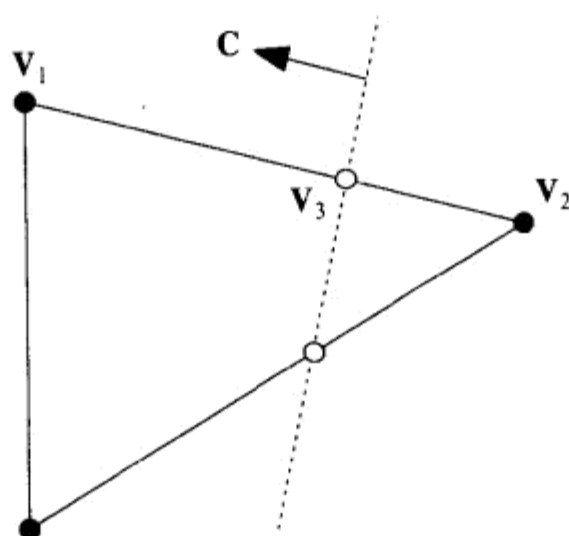


图 4.8.2 沿着平面剪裁多边形

假设顶点  $V_1$  位于剪裁平面  $C$  的正面， $V_2$  位于剪裁平面  $C$  的背面。连接点  $V_1$  和  $V_2$  的线段上的一点  $V_3$  可以表示为：

$$V_3 = V_1 + t(V_2 - V_1) \quad (4.8.5)$$

其中参数  $t$  满足  $0 \leq t \leq 1$ 。我们想知道  $t$  取何值时，点  $V_3$  位于平面  $C$  上。如果我们将  $V_i$  作为  $w$  坐标为 1 的齐次向量，则只需要找到  $t$  值，使得四维点积  $C \cdot V_3$  等于 0。将  $V_3$  带入公式 4.8.5 解得  $t$  为：

$$t = \frac{C \cdot V_1}{C \cdot (V_1 - V_2)} \quad (4.8.6)$$

（注意：差  $V_1 - V_2$  有一个  $w$  坐标为 0。）将  $t$  的值代回公式 4.8.5 得到一个新顶点。

### 4.8.3 实现代码



附带光盘上的源代码演示了本文描述的算法，它通过一个叫作 Decal 的 C++ 类来实现。这个类的构造函数以贴花中心  $P$ 、法线方向  $N$  以及切线方向  $T$  作为参数，还使用了贴花的宽、高、深度。用公式 4.8.1 和 4.8.2 计算边界平面以后，构造函数剪裁所有可能受影响的表面，并将结果网格保存在一个新三角形数组中。顶点颜色和纹理映射坐标由公式 4.8.3 和 4.8.4 指定。源代码生成的焦痕贴花如图 4.8.3 所示。

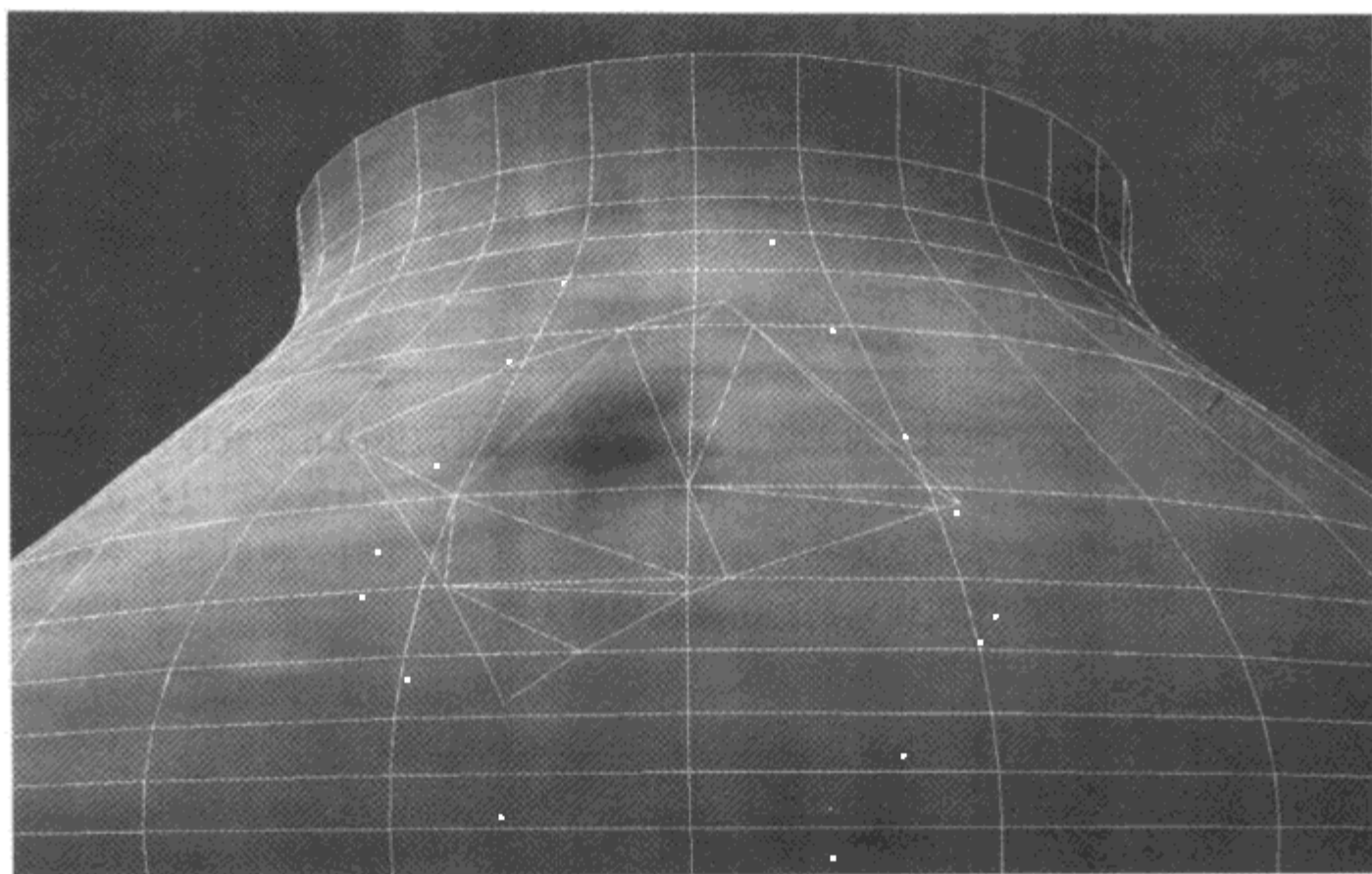


图 4.8.3 曲面的焦痕贴花

### 4.8.4 参考文献

[Lengyel00] Lengyel, Eric, "Tweaking a Vertex's Projected Depth Value," *Game Programming Gems*, Charles River Media, 2000, pp. 361~365.



## 4.9 用天空包围盒渲染远景

Jason Shankel, Maxis  
shankel@pobox.com

**请**想像一个发生在西藏高原之巅的游戏场景。在远方，我们可以看到绵延的喜马拉雅山脉，天空中白云飘飘，山谷里村落点点。或者请想像一个发生在宇宙深处的游戏，天鹰和猎户星云在数以光年计的距离以外闪闪发光。如此美妙场景可以增加 3D 游戏的美感和身临其境的感受。

渲染远距离 3D 场景可以用天空包围盒 (skybox) 来完成。本文解释了天空包围盒技术的原理，并描述了渲染天空包围盒场景的几种方法。

### 4.9.1 基本技术

天空包围盒的思想非常简单。将远景渲染成 6 种纹理，每种纹理应用到立方体的一面。摄像机置于立方体的中心。摄像机可以在立方体内任意旋转，但是不能离开立方体中心。当渲染场景时，投射到天空包围盒各个面上的图像表现远景，正像投影到天文馆天花板上的图像表现出夜空一样 (如图 4.9.1 所示)。

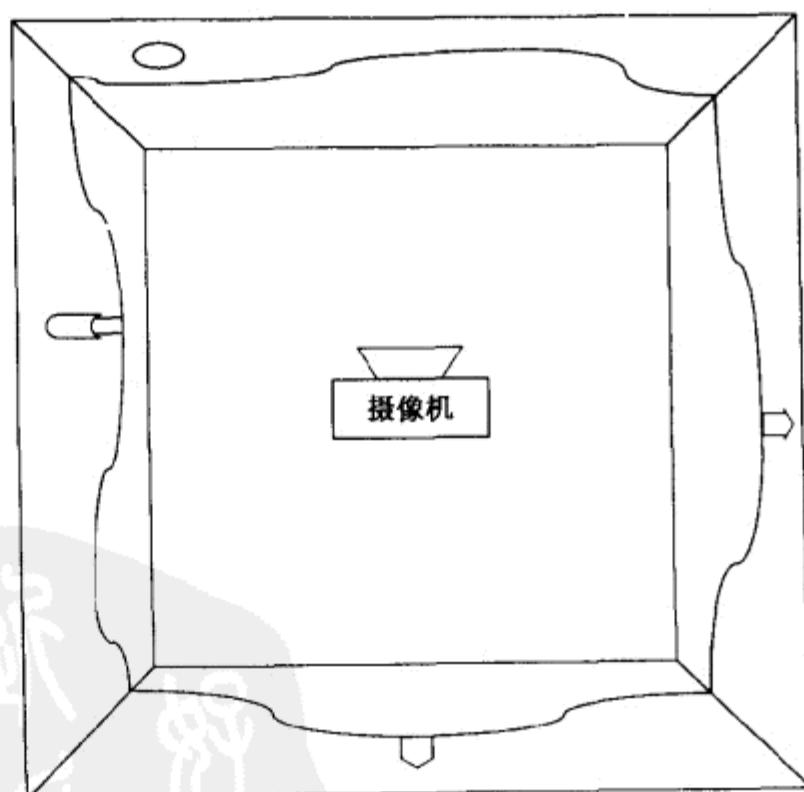


图 4.9.1 天空包围盒的俯视图：在天空包围盒各侧面上渲染远景地形，摄像机置于中心位置



### 4.9.2 天空包围盒分辨率

理想状态下, 天空包围盒 (skybox) 物体中的一个贴图像素 (texel, 即纹理元素) 应当映射到屏幕上的一个像素。下式可用来决定已知屏幕分辨率下的理想包围分辨率:

$$\text{skybox Res} = \frac{\text{screen Res}}{\tan(\text{fov}/2)}$$

其中  $\text{skyboxRes}$  是以贴图像素为单位的天空包围盒侧面的分辨率,  $\text{screenRes}$  是以像素为单位的屏幕宽度,  $\text{fov}$  是水平视场角度。例如, 一个运行在  $640 \times 480$  的分辨率下、 $90^\circ$  视场中的游戏, 理想天空包围盒分辨率为 640 贴图像素。

有些 3D 系统将纹理分辨率限制为  $256 \times 256$  贴图像素。这使得包围的纹理渲染到屏幕上时, 可能会被显著拉伸。在有些应用场合, 这样的拉伸可以接受。在其他场合, 可能希望分割天空包围盒的表面以提高纹理分辨率 (如图 4.9.2 所示)。

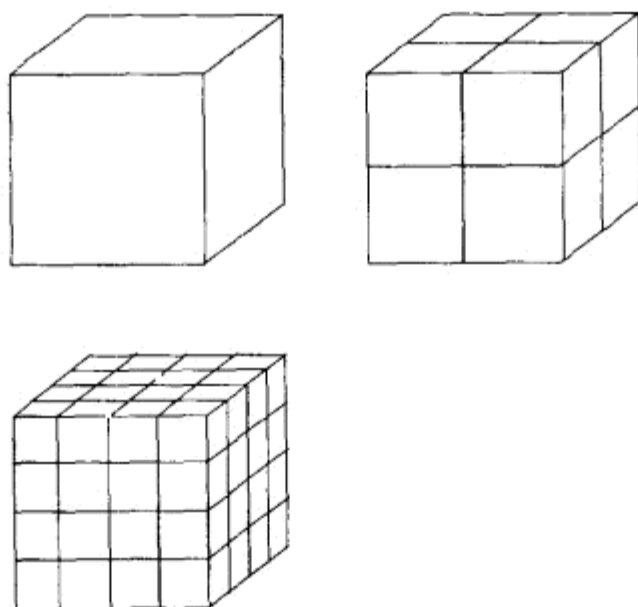


图 4.9.2 分割天空包围盒的表面将以增加纹理内存和多边形数量为代价来改善图像质量。

左上图: 每边 1 个纹理。右上图: 每边 4 个纹理。左下图: 每边 16 个纹理

### 4.9.3 天空包围盒大小

因为摄像机总是位于中心, 所以不管采用多大的天空包围盒, 只要它落在视锥内即可。当天空包围盒增大时, 它的侧面成比例地向远处延伸, 而它在画面中所占部分保持不变。

选择合理的天空包围盒大小, 使得最远处的点及包围盒的角比摄像机远剪裁面更近。下式可用来计算包围的最大尺寸:

$$\text{skyboxWidth} \leq \frac{2\sqrt{3}}{3} z_{\text{far}}$$

其中  $\text{skyboxWidth}$  是世界坐标系为单位的天空包围盒一个边的长度, 而  $z_{\text{far}}$  是到世界坐标系中远处错切平面的距离。

#### 4.9.4 渲染场景

绘制天空包围盒的最简单技术是渲染一个具有纹理的普通立方体，与世界坐标轴对齐，并以摄像机所处的世界位置作为中心。首先渲染天空包围盒。不需要先清空颜色缓冲，因为天空包围盒总是充满整个屏幕。

当渲染天空包围盒时，要禁止深度测试和深度写入，因为天空包围盒几何体与场景中的其他几何体不会完全相交。光照、雾化和其他类似效果也要禁止。渲染天空包围盒图像时不进行任何大气效果调整。

前面提到，一面一纹理的简单立方体可能会产生明显的纹理拉伸。可以采用纹理过滤来减弱纹理拉伸效应。纹理过滤通过对相邻的贴像素进行采样并将其混合在一起，减少拉伸图像的锯齿。

纹理过滤可能导致沿着天空包围盒的边缘（该处是两种纹理交叉的地方）产生拼合裂缝。这些裂缝由纹理的不合适卷绕产生。纹理叠加在纹理边缘上确定哪些贴图像素为“邻近”贴图像素。如果纹理叠加设置重复（OpenGL 中的 `GL_REPEAT`），则位于一条边上的贴图像素将与相对边上的贴图像素混合在一起，导致纹理之间凸现拼接裂缝。

如果可能的话，在 OpenGL 中将纹理卷绕设成 `GL_CLAMP_TO_EDGE_EXT`。这将固定对边的过滤，消除来自边界贴图像素的影响。`GL_EXT_texture_edge_clamp` 模式是 OpenGL 的一个扩展，所以不是在所有系统上都可以使用。如果 `GL_CLAMP_TO_EDGE_EXT` 不可用，纹理过滤应该设成 `GL_CLAMP`，它通过贴图像素的恒定边界颜色融合边上的贴图像素。

#### 4.9.5 立方体环境映射

立方体环境映射（cube environment mapping）是传统天空包围盒渲染技术之外的另一种方法。立方体环境映射将天空包围盒的 6 个面融合成一张贴图。

通常而言，纹理坐标被定义为二维纹理映射中的偏移，(0, 0)点是纹理的一个角，(1, 1)是对角。在立方体环境映射中，6 种纹理融合起来形成一个立方体，纹理坐标定义为三维向量，从立方体的中心指向其表面的一点（如图 4.9.3 所示）。

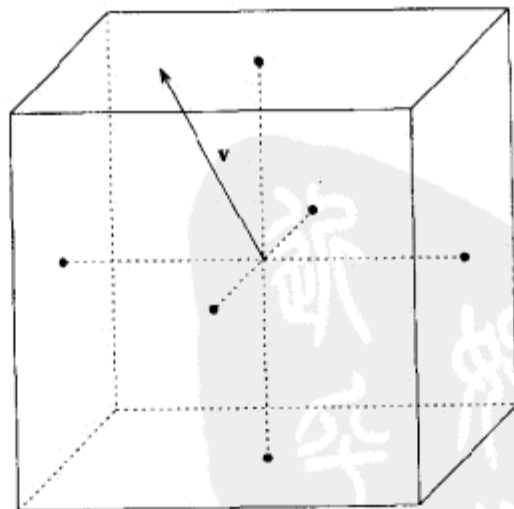


图 4.9.3 立方体环境映射的纹理坐标定义为以原点为中心、指向立方体表面的向量

立方体环境映射可用来渲染远景，并且将天空包围盒的折射光投影到附近闪光物体上。

立方体环境映射限定天空包围盒的每一个表面为一个 2D 纹理，因此分割表面以增加分辨率是不可能的。好在大部分立方体环境映射系统支持  $256 \times 256$  以上的纹理分辨率。

#### 4.9.6 生成天空包围盒纹理

---

一些软件包提供对生成天空包围盒纹理的直接支持。对那些不支持的软件包，手工生成天空包围盒图像也很简单。

输出图像的大小应该设置为纹理分辨率（例如  $256 \times 256$ ），并且摄像机的视场设置为  $90^\circ$ 。生成 6 个渲染图像，用摄像机沿三基轴的 6 个方向（左、右、上、下、前、后）进行观察。

#### 4.9.7 结论

---

在大部分 3D 游戏中，场景的作用至关重要。渲染远景不仅增加了视觉的美感，也为玩家提供了一种在虚拟环境中更加完善的方向性。天空包围盒提供了一种渲染远景的既经济又有效的方法。

#### 4.9.8 源代码

---



提供的源代码使用了 OpenGL 和 GLUT。示例程序对 2D 和立方体环境映射渲染均作了演示。

OpenGL 通过 `GL_ARB_texture_cube_map` 扩展支持立方体环境映射。示例程序在开始时检查这个扩展是否存在。在立方体环境映射模式下，将闪光物体放在摄像机的视场中，来演示折射效果。

通过 Terragen（Copyright © 1997–2000 Planetside Software）生成天空包围盒纹理。



## 4.10 自阴影角色

---

Alex Vlachos, David Gosselin, Jason L. Mitchell; ATI Research  
alex@vlachos.com, gosselin@ati.com, jasonm@ati.com

许多实时游戏缺乏渲染自阴影角色这一重要视觉效果。本文讲解投影纹理方法，可以实时处理游戏角色凸的部分的自阴影。将角色分解成凸的子部后，以光源为视角，用不同的  $\alpha$  值将子部渲染成纹理。当从玩家的视角渲染角色时，再把该纹理投影回角色凸的子部，模拟每个子部相对于光源的遮挡。由于是基于映射，而不是基于模板体积，无须额外工作即可与表面细分技术兼容。

### 4.10.1 研究回顾

---

阴影的投影技术始于 Williams ([Williams78])，并存在许多对之进行改进的论文。本文摒弃了基于深度的阴影映射技术，而是使用一种类似于优先缓冲 (priority buffer) 的技术 [Hourcade85]。文献 [Lengyel00] 也使用了类似的分割技术。

### 4.10.2 角色几何分割

---

算法的第一步是对角色几何体进行分组。美工标记骨骼形成几何体组，该几何体组将遮住其余几何体组。必须精心选择这些组，以减少模型的分割数，但是要保留遮挡关系。在预处理过程中，以最大重量的骨骼作为每个多边形的主骨骼。在标记主骨骼后，将多边形放在合适的组内。图 4.10.1 是一个分割模型的示例。

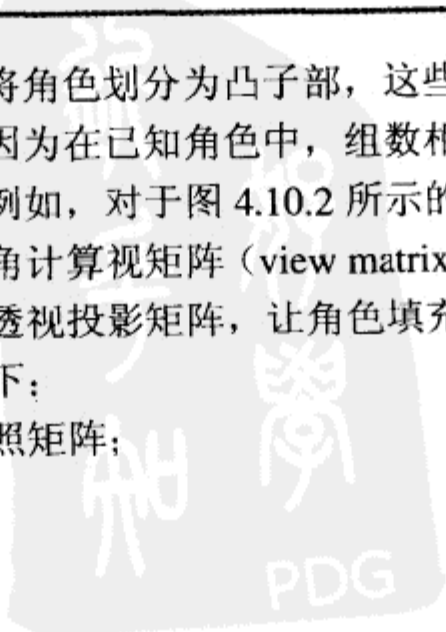
### 4.10.3 渲染纹理

---

一旦按逻辑将角色划分为凸子部，这些组将依照其与光源的关系从前到后进行排序。因为在已知角色中，组数相当少。所以，这种排序的计算负担并不很重。例如，对于图 4.10.2 所示的角色，仅仅要对 6 个组排序。然后从光源的视角计算视矩阵 (view matrix)，它的中心放在物体包围盒中心。同时也生成透视投影矩阵，让角色填充尽可能多的可渲染纹理。

基本步骤如下：

- (1) 计算光照矩阵；



(2) 用(255, 255, 255, 255)清空可渲染纹理;

(3) 将视点放在光源上, 从前向后 alpha 值以 0, 1, 2, 3, 4, …… 254 增加 (255 作为 alpha 清色值), 以此画出每个凸的子部表面的 alpha 纹理。

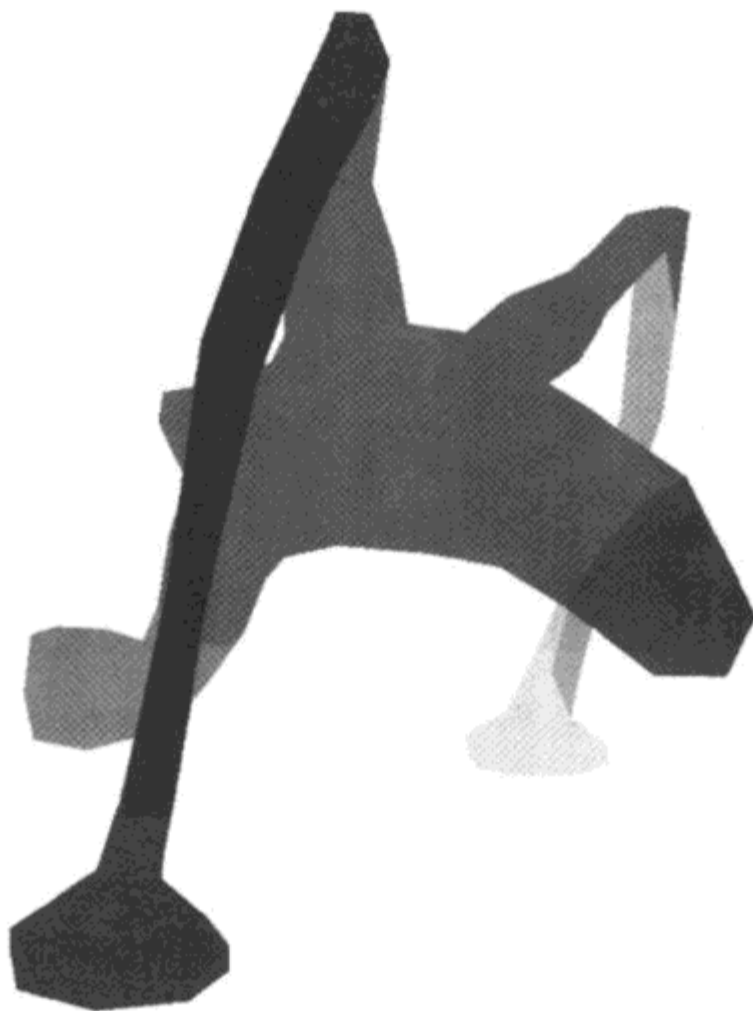


图 4.10.1 分割模型

注意, alpha 值的可以递增覆盖整个有效 alpha 值范围, 以对应需要渲染的组数。这将有助于消除渲染的不自然性, 下文将对此作出解释。

#### 4.10.4 渲染角色

一旦创建了阴影贴图, 下一步就是要使用阴影贴图将角色渲染到帧缓冲中。对于蛮力 (brute-force) 方法, 可用三步来完成。第一步将完全光照角色渲染到帧缓冲。第二步对每一片的 alpha 值进行检测, 如果 alpha 值小于阴影缓存中的值, 则绘制该片。这样做基于两个理由。第一个理由是为了只使用对物体起阴影作用的阴影缓冲。因为阴影缓冲按不同的 alpha 值分段<sup>1</sup>, 所以把 alpha 比较模式设置为“小于”, 就可以得到之前所有物体组的阴影 (这些物体组是更靠近光源的物体)。使用 alpha 测试的第二个原因是为了有效地把阴影缩小半个像素<sup>2</sup>, 从而避免由于双线性过滤贴图而产生的不必要的 alpha 值。

使用渲染阴影贴图时设置的灯光视角矩阵 (view matrix) 和投影矩阵 (projection matrix),

1. 注: alpha 值越大的阴影缓冲值代表物体像素离灯光距离越远。

2. 注: 当使用阴影缓冲贴图绘制物体时, 因为硬件有双线性过滤功能, 显卡会在几个 alpha 贴图像素之间做线性插值, 得到最终绘制的屏幕像素。这样会造成阴影边界上的绘制误差。

来转换物体的世界坐标位置矩阵。得到的结果用来设置（阴影）贴图矩阵。通过 alpha 检测的像素使用黑色绘制，用以有效地屏蔽在第一遍渲染中绘制的带有光照的物体像素。最后一遍渲染使用角色的基本（贴图）颜色与环境光照调制，用以加亮物体阴影的区域。图 4.10.2 显示了使用此方法渲染的动画中的几帧。这种方法显然是一种蛮力法实现。它可以使用像素着色器（pixel shader）来进行优化。

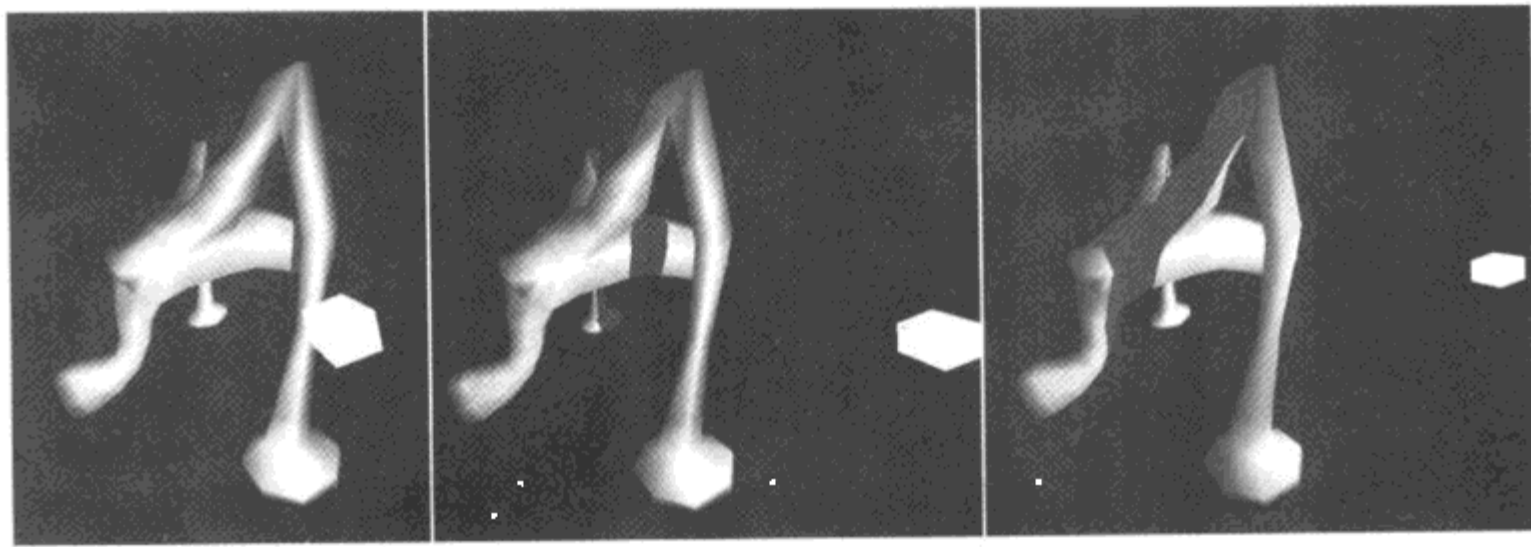


图 4.10.2 从 3 个不同角度渲染带有阴影的造型

#### 4.10.5 结论

---

我们描述了一种采用投影纹理技术，渲染具有自阴影的动画造型。通过将造型分割成凸起细部，避免在纯基于深度方法中出现的混淆部分。这种方法也非常通用，可应用于任何支持纹理渲染和投影纹理映射的图形硬件。

#### 4.10.6 参考文献

---

[Hourcade85] Hourcade, J. C., and Nicolas, A., "Algorithms for Antialiased Cast Shadows," *Computers and Graphics*, vol. 9, no. 3, pp. 259~264.

[Lengyel00] Lengyel, Jerome E., "Splatting for Soft-Edged Shadows," Microsoft Technical Report, 2000.

[Williams78] Williams, Lance, "Casting Curved Shadows on Curved Surfaces," *Computer Graphics (SIGGRAPH 78 Proceedings)*, vol. 12, no. 3, 1978, pp. 270~274.





## 4.11 经典的 *Super Mario 64* 游戏第三人称控制和动画

Steve Rabin, Nintendo of America

steve\_rabin@hotmail.com

经典的 *Super Mario 64* (超级玛丽) 游戏控制方案轻而易举地成为从第三人称视角控制 3D 人物角色的一种最直观方式。尽管 *Super Mario 64* 游戏并没有发明这项技术, 但它确实在尽力去改善它, 并在众多的玩家中普及这项技术。它经常被用作一种极具可玩性控制方案的测试尺度。它带给首次游戏者的那种美妙感觉, 可不是一个小成就。

本文将讲解从第三人称视角来控制动画游戏造型的基本问题。尽管看起来这十分简单 (只是照搬 *Super Mario 64* 而已), 但事实上并不像想像的那么简单。这其中有很多智慧的结晶, 需要数周的试验和 排错才能挖掘出来。

### 4.11.1 设置

经典的 *Super Mario 64* 控制中包括一个自由浮动的摄像机, 不管游戏角色出现在哪儿, 摄像机都会跟随在其周围。可以设想成, 这个摄像机悬浮着, 用一根弹簧索系在角色身上。随着角色从站立姿势变为跑步姿势, 摄像机也逐渐加速, 尽可能与角色保持距离恒定。当角色转向的时候, 摄像机也尽可能保持位于其后, 慢慢做圆周运动。

这个控制方案的核心是: 角色将按照我们指定的任意路径行走, 就好像我们坐在摄像机内部一样。因此, 如果我们按控制器的向上键, 角色将离开摄像机。如果按向右键, 角色就会从摄像机的视角向右移动。棘手的问题是, 摄像机总在行动, 导向的摄像机的控制需要映射到角色实际活动的世界坐标系中。

### 4.11.2 转换控制器的输入

基本的编程问题包括把玩家的输入 (上、下、右、左) 从摄像机的视角转换到自然方位。为了做到这一步, 我们分别需要与摄像机的前进方向和右进方向对应的向量。

由于控制器任何向上或向下的运动都直接映射为沿摄像机的前进向量的转换, 这些向量就变得非常有用。同样, 控制器任何向右或向左的移动也映射为沿摄像机右向向量的转换。不过, 由于一般角色总是在一个平面上

移动，这些摄像机向量就不应该包含任何高度的信息（一般用 y 轴来表示）。图 4.11.1 演示如何相对于摄像机和角色考虑相应的控制器输入。

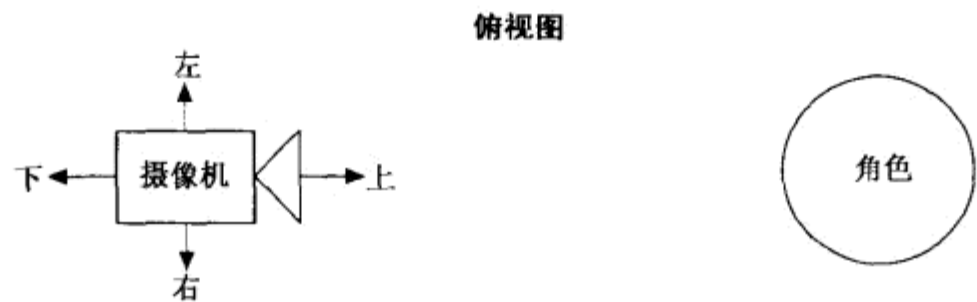


图 4.11.1 控制器与角色之间针对摄像机的关系

由于摄像机的前向向量就是摄像机面对的方向，所以非常容易找到它。一旦找到了前向向量，我们就可以标出垂直轴的 0 坐标，并且重新标准化（normalize）该坐标轴。可以根据如下关系式确定摄像机的右向向量：

```
right.x = forward.z
right.z = -forward.x
right.y = 0.0
```

由于控制器输入传统上以数字的形式提供，有必要简要解释一下。可以认为用户在模拟杆上的输入在各根轴上的基本变化范围是-1.0~1.0。实际上，大多输入方案采用-128~127 或-32768~32767 的整数。但是，针对我们的特殊目的，需要将它们转换成-1.0~1.0 的浮点数。然而，我们必须知道，向上和向右的都是正数。控制器的布局如图 4.11.2 所示。

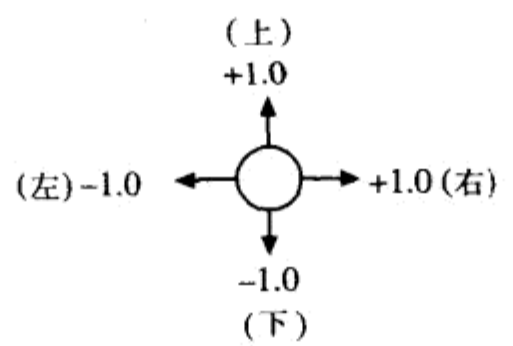
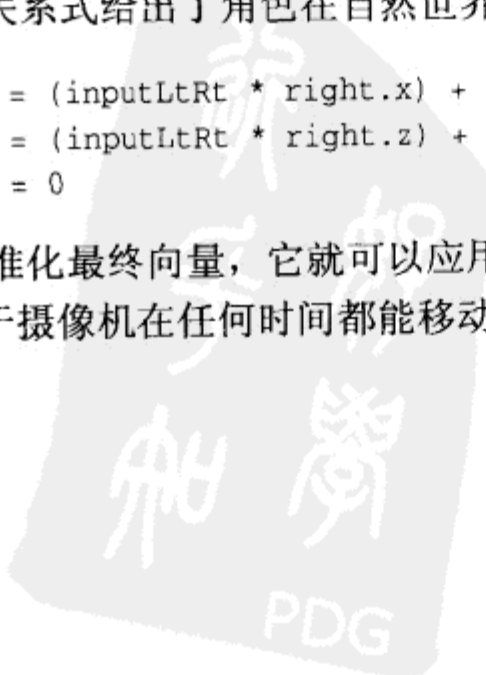


图 4.11.2 模拟控制器的值

有了摄像机的前向和右向向量，就能通过重新映射控制器的输入得到最后的自然坐标向量。下面的关系式给出了角色在自然世界坐标体系中的移动方向。

```
final.x = (inputLtRt * right.x) + (inputUpDn * forward.x)
final.z = (inputLtRt * right.z) + (inputUpDn * forward.z)
final.y = 0
```

一旦标准化最终向量，它就可以应用于角色上，从而保证它按正确的方向移动。我们应该记住，由于摄像机在任何时间都能移动，每一帧都必须计算一个新的前向向量和右向向量。



### 4.11.3 旋转角色

任一时刻的用户输入均给出一个期望的运动方向。量取一段时间中角色的运动距离，可以确定角色的速度。移动角色的一个简单的模型是分别保持方向与速度。这有助于把它们分开，而不是把两者表示为一个向量，所以零速度仍然具有方向性。显然，如果我们试图用零向量来引导角色，就会发生问题。

经过一段时间对输入取样，角色的方向应该集中在用户期望的方向上。比如，如果角色正面向北，而玩家把他转向东，角色应该利用某种衰减转向东。可以基于帧频，在当前方向上增加一定比例的期望方向，会获取非常好的衰减效果。在图 4.11.3 中，请注意角色是如何在开始时快速旋转，而后在一个短暂的时间内，慢慢衰减到期望的方向上的。

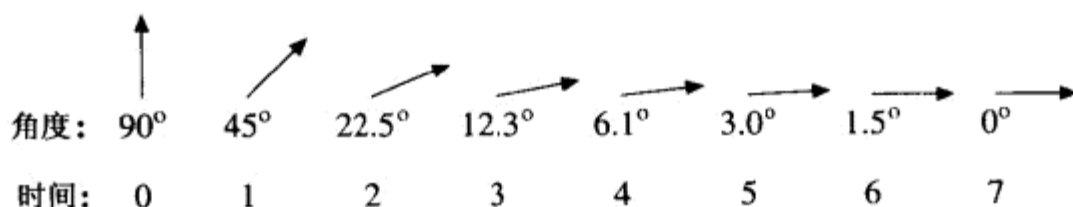


图 4.11.3 角色通过指数衰减随时间的旋转

在控制角色时，图 4.11.3 中的柔性衰减导致平滑的响应感觉。优化收敛于目标方向的速率可以获得所期望的感觉。图 4.11.3 中的衰减函数可以表示为如下公式：

```
new_dir = (old_dir * 0.5) + (desired_dir * 0.5)
Normalize( new_dir )
```

通过考虑帧频，可以让衰减函数与帧频变化相当一致。不幸的是，我们还不能正确处理帧频 100% 上下波动的情况。下面是一个与帧频变化呈半弹性的衰减方程（在特定窗口下）。

```
new_ratio = 0.9 / FPS
old_ratio = 1.0 - new_ratio
new_dir = (old_dir * old_ratio) + (desired_dir * new_ratio)
Normalize( new_dir )
```

不管帧频是多少，这个新衰减公式会在一秒后把角色朝期望的方向旋转 90°。百分比需要优化，但是不论帧频是多少，我们都把百分比调整到帧速率窗口的中点。但是，当帧频与那个中点差别很大时，“感觉”就会慢慢失真，就会变得过分敏感或过分迟钝。如果帧频变化太大，就可以考虑用图形更新来去耦计算，以保持一个平稳的速率。

当要求角色从当前方向转动 180° 时，平滑旋转并不总是最好的方案。由于要旋转一个如此大的角度，将导致一种单调、不自然的持续旋转。最好利用一个过渡动画来完成比较大的方向改变。在 *Super Mario 64* 中，当 Mario 全速奔跑时被命令反向，即发生这种大角度旋转，会导致一个滑停动作，随后是一个 180° 原地转向。但是，在这个过程中，再次移动前存在一个瞬间停顿的速度减慢问题。

到目前为止，我们所描述的旋转方法相当敏感，容易受帧频的影响，尽管我们已尽可能地减少这种影响。比如，如果您记录一个玩家的输入，同时以稍微不同的帧频反馈到游戏，

结果也会稍有差别。如果您仅用输入来记录游戏，而且希望从不同的摄影机角度来重玩游戏（也许导致帧频变化），结果是很不同的。解决这个问题的方案是，要么锁定模拟频率，要么使用高度控制的旋转方法，比如，设定一个特定的角速度。但是，这意味着为了追求更高的可预测性，而放弃一些平滑的加速度效果。

#### 4.11.4 角色移动

当角色平滑转动的时候，我们需要真正地移动他。要沿着角色目前的方向，用当前速度向前移动他。角色的平滑旋转基本上是以适当的方向围绕臀部进行，同时用当前速度驱动他向前移动。

与控制旋转时使用的衰减函数一样，也需要某种衰减函数来控制速度。但是，与旋转不一样，需要把速度保持在最高速。需要降低当前速度进行  $180^\circ$  和一般的转动。为了使控制反应灵敏，衰减函数应该很快地收敛到期望速度（快于  $1/10\text{ s}$ ）。

响应度是一个很有意思的量。在理想状态下，角色应该即时地对玩家输入做出反应。不过，这个要求不仅不现实而且强人所难。增加一个指数衰减，与实际过渡动作略有不同，但保持反应灵敏度。请注意，这与物理运动不同。角色可以强有力地控制他们的运动，但实际上不会像汽车和宇宙飞船那样加速和减速。图 4.11.4 演示了控制器输入下角色的速度。注意响应度等同于快速地增加到期望速度，然后平滑地达到该速度值的过程。

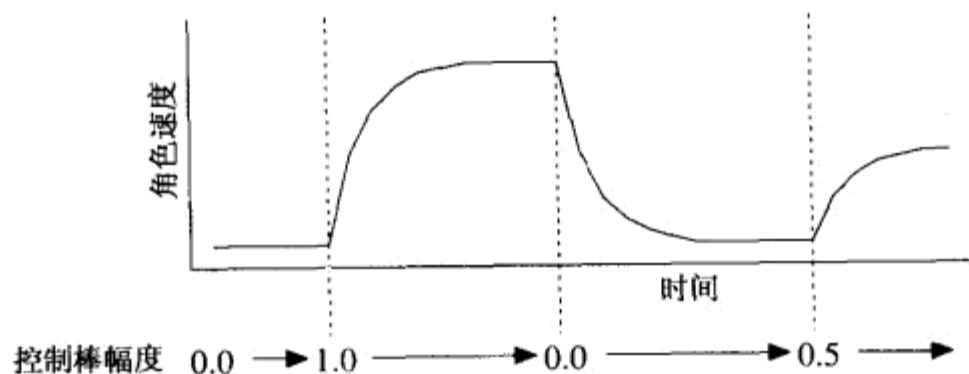


图 4.11.4 基于指数衰减控制器输入的角色速度

#### 4.11.5 角色动画

最简单的动画方案是有两个动画循环：站立和奔跑。当角色的速度是零的时候，播放站立动画。当角色的速度大于零时，播放奔跑动画。当然大多角色的动画应该更加精细，但这是良好的起点。

下面是一些让动画看起来更加合理的技巧：

- (1) 把动画的速度（脚的移动）与角色的最高速度相匹配。当角色移动的时候，大多数时间以最高速度进行，所以该循环应该在这个速度上看起来最佳。在这个速度上，他们的脚应该接触地面，看起来与地面粘住，直到他们再抬起脚。
- (2) 可以改变动画的速度匹配角色的速度，但是很多动画只是在很小的速度范围内看起来是合理的。速度变化范围很小，而对每种运动并不是在  $0\text{ m/s} \sim 8\text{ m/s}$  的整个范围里

变化。奔跑时的速度一般是  $5\text{m/s} \sim 8\text{m/s}$ ，小跑一般是  $3\text{m/s} \sim 5\text{m/s}$ ，行走时一般是  $1\text{m/s} \sim 3\text{m/s}$ 。真正的范围只能经验地根据实际动画数据确定，因而这一范围具有极大的主观性。

- (3) 当角色突然高速转向时，角色应该从算法上让他倾斜到转向方。这有助于强调速度，而且看起来也更加自然。另一个增强转向效果的方法是，从算法上让角色将头稍稍朝向转弯方向。由于角色一般会预见并朝向下一个位置，这种方法进一步增加了真实感。
- (4) 实现动画过渡（站立到奔跑，行走到奔跑）需要一定技巧。通常他们都会妨碍响应度。一个简单明确的解决方案是在过渡时从算法上构建一个中间帧。在不影响视觉外观（Façade）和感觉的情况下消除了边缘过渡。计算中间帧是一个众所周知的问题，很多动画编程书籍中都会提到。
- (5) 通过正确设计行走或奔跑的循环，实际上可以实现从站立到行走或到奔跑的平滑过渡。图 4.11.5 显示了循环中第一帧的样子。请注意当右脚差不多抬起时，左脚是如何立在地面上的。这个特定的关键帧使得从站立到行走或到奔跑时，动画的突变最小。但是，因为移动的角色可能在任何时间停止，这无助于平滑过渡到站立。快速把角色插值到站立的位置上，可以过渡到站立，同时小心不要让脚穿透地面。

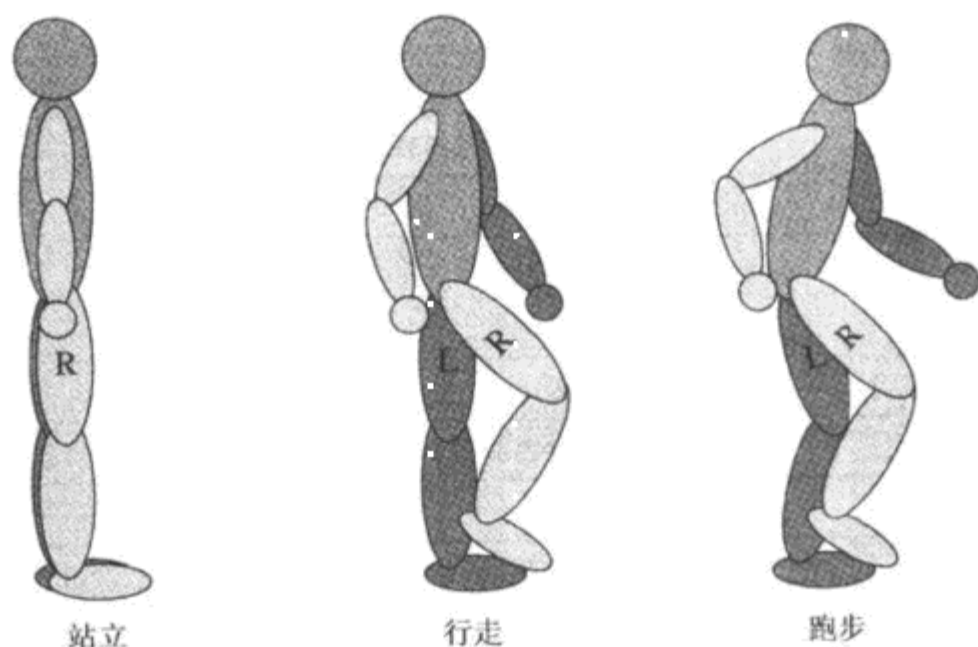


图 4.11.5 平滑过渡的第一帧

- (6) 可以设计一个更加精细的过渡系统，在这个系统里，存在一些真正的过渡动画，比如从站立过渡到奔跑、从奔跑过渡到站立等。当玩家命令角色开始移动的时候，播放从站立到奔跑的动画，随后是奔跑的循环。当玩家命令角色停止的时候，角色必须平滑过渡到从奔跑过渡到站立的动画，而不管奔跑循环的当前帧是什么。解决这个简单方法是，拥有两个不同的奔跑过渡到站立的动画，每个动画中的起始前脚不同。在这种方法下，从奔跑过渡到站立动画前，只需要播放奔跑循环的一半。线性过渡需要时间，但如果强迫动画播放更长的时间，就会牺牲一定的灵敏度。而且，过渡动画必须非常快，否则，灵敏度就会受到进一步的影响。根据期望得到的游戏效果和感

觉，应该通过经验和算法上的过渡进行试验获得恰当的融合点。

#### 4.11.6 *Super Mario 64* 动画分析

表 4.11.1 列出了 *Super Mario 64* 中的一些动画。这些信息是在不具备任何代码知识的情况下通过观察游戏得到的。当然，为了解析控制和动画技术，可以按这种方式分析任意游戏。

**表 4.11.1** *Super Mario 64* 中的一些动画

速度	动画
站立	若干静止动画
中速	小跑
快速	前倾式跑步
极快	明显前倾式跑步
过渡	动画
站立到脚尖着地	一个站立循环，接着是脚尖着地
站立到小跑	快进到小跑
站立到跑步	快进到跑步，且冒出一股烟
脚尖着地到站立	快进到站立
小跑到站立	快进到站立
跑步到站立	快进到滑停，然后逐渐静止
脚尖着地到小跑	快进到小跑
小跑到跑步	快进到跑步（匹配脚步循环），且冒出一股烟
小跑到脚尖着地	快进到脚尖着地
跑步到脚尖着地	快进到脚尖着地
跑步到小跑	快进到小跑（匹配脚步循环）
动作	引起行为
站立时的任意转向	快进到目标方向
脚尖着地时 180°转弯	平滑 U 形转弯
小跑时 180°转弯	平滑 U 形转弯
跑步时 180°转弯	刹车，并向后转
脚尖着地时小于 180°的转弯	平滑旋转
小跑时小于 180°的转弯	平滑旋转
跑步时小于 180°的转弯	平滑旋转
跑步时转弯	倾斜转弯



#### 4.11.7 结论

---

毫无疑问,甚至在 *Super Mario 64* 出现以前,很多优秀的游戏就已经实现了经典的 *Super Mario 64* 第三人称控制。但是,在试图编写这样的代码以前,找一个例子来模仿,可以很快就熟悉所使用的方法。尽量辨别不同控制器输入下所发生现象的细微差异。要注意过渡、移动/旋转的加速/减速,行走/小跑/奔跑的速度差异、缓转弯和急转弯以及完全松开控制器后的衰减等。有一个优秀的例子来模仿非常重要,因为你可以有一个稳定的对象临摹,并最终改进它。

#### 4.11.8 参考文献

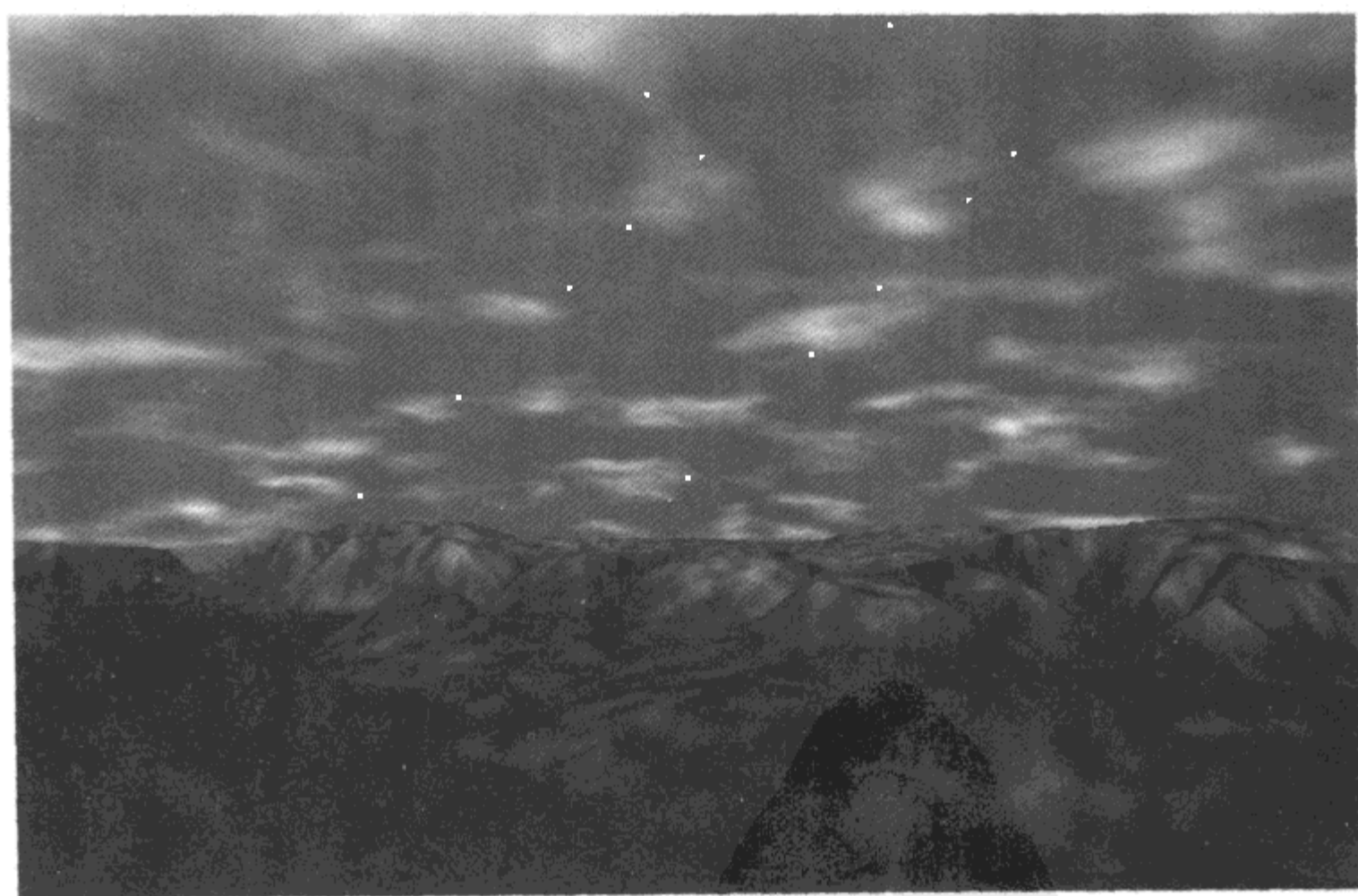
---

[Maestri99] Maestri, George, *Digital Character Animation 2: Essential Techniques*, New Riders Publishing, July 1999.





## 图形显示





## 绪 论

---

D. Sim Dietrich Jr., Nvidia 公司

[sdietrich@nvidia.com](mailto:sdietrich@nvidia.com)

本章（图形显示）主要是为读者提供实用而有创意的实时图形技术。另外，尽量提供足够多的信息，以便读者能理解这些技术是如何实现的，以及怎样扩展、更改这些技术，以满足读者特定的需求。进一步的目标是，介绍在未来 5 年左右仍适用的技术。能否如愿以偿，让时间来检验一切，不过我抱谨慎的乐观态度。

“图形显示”这部分内容包括实时三维图形技术的几个不同方面，但许多文章都归结到一个特点的技术——渲染纹理（render-to-texture）功能。这并不是有意的偏爱，而是由于这项技术本身灵活有用。据此，我相信，渲染纹理是未来许多实时图形所基于的基本运算之一，因为它挖掘出了现代图形硬件难以置信的实时性威力，所产生的效果在以前是不可能达到的。

我与其他作者一样，为有幸替“游戏编程精粹”丛书撰稿而表示感谢，同时，我们也希望能为实时图形在黄金时代的技术发展贡献微薄之力。



## 5.1 卡通渲染：实时轮廓边缘检测与渲染

Carl S. Marshall, 英特尔架构实验室

carl.s.marshall@intel.com

**轮廓检测与渲染** (silhouette detection and rendering) 是为 3D 卡通渲染添加特定风格外观 (stylized look) 的关键部分。轮廓边缘检测 (silhouette edge detection) 的基本概念是查找对于描述模型轮廓比较重要的边。许多卡通动画师通过跟踪模型周边的黑色外形来显示这些轮廓。这一卡通渲染核心描述了若干边缘检测技术：基于边的检测方法、可编程顶点着色技术以及高级纹理技术。每种技术各有优缺点。

### 5.1.1 着墨器 (Inker)

术语着墨 (inking) 和上色 (painting) 源自于传统的卡通创建过程。着墨指为场景中的人和物描出轮廓的过程，而上色指在轮廓内部上色或填色的过程。这种技术核心的重点是卡通渲染的着墨过程。Adam Lake 的文章《使用纹理映射的卡通渲染和可编程顶点着色器》讲述了上色过程。这两种技术一起组成了一个完整的卡通渲染引擎。图 5.1.1 演示了一个 3D 鸭子模型的着墨和上色过程。

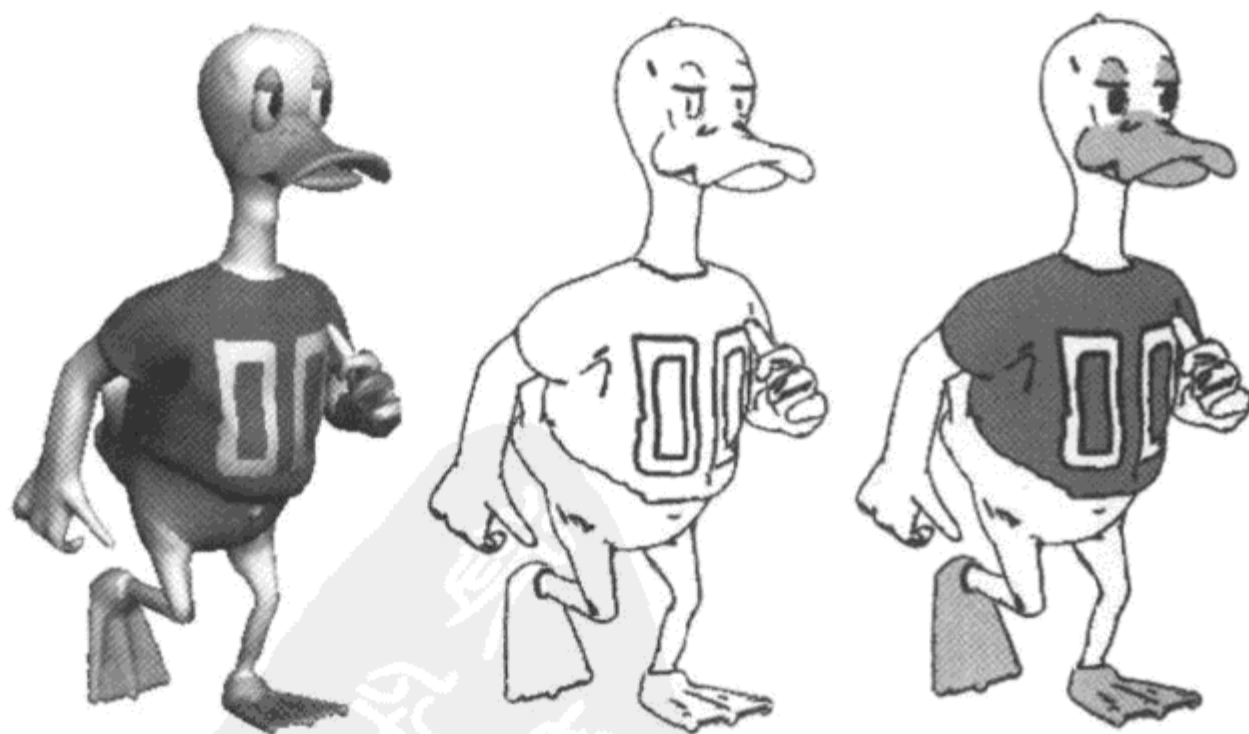


图 5.1.1 Gouraud 着色 (Gouraud shading) 的鸭子 (左图)、脸部着色为背景的已着墨鸭子 (中间图) 以及使用上色器 Flat 着色的鸭子 (右图)



一般的着墨 (inking) 过程包含两部分。第一部分为轮廓检测, 第二部分是轮廓渲染。该核心技术使用若干可以对多样式以及非多样式网格实时执行的技术描述着墨过程。

### 5.1.2 重要的边

轮廓边缘检测 (Silhouette Edge Detection, SED) 是着墨器的主要组成。与轮廓的检测与渲染一起, 画家还要突出一些重要的模型边: 折缝边 (crease edge) 和边界边 (boundary edge)。如绪论里所描述的, 轮廓边缘是形成模型外形的边, 但也可以包含一些内部边, 如图 5.1.2 所示。改变轮廓检测的一个重要方面是, 视角是相关的。这意味着, 随着照相机或模型的运动, 轮廓必须针对每次改变重新检测与渲染。折缝角边指相邻两个面之间的角度位于已知阈值之内的边。它称做二面角 (dihedral angle)。使用用户自定义阈值来检测折缝角。边界边指不连续面接壤处的边界。当网格中两个三角形接壤处的边共用相同的 2 个顶点, 但是两个顶点的贴图坐标、材质或者法线不一样, 这两个三角形就是不连续面。

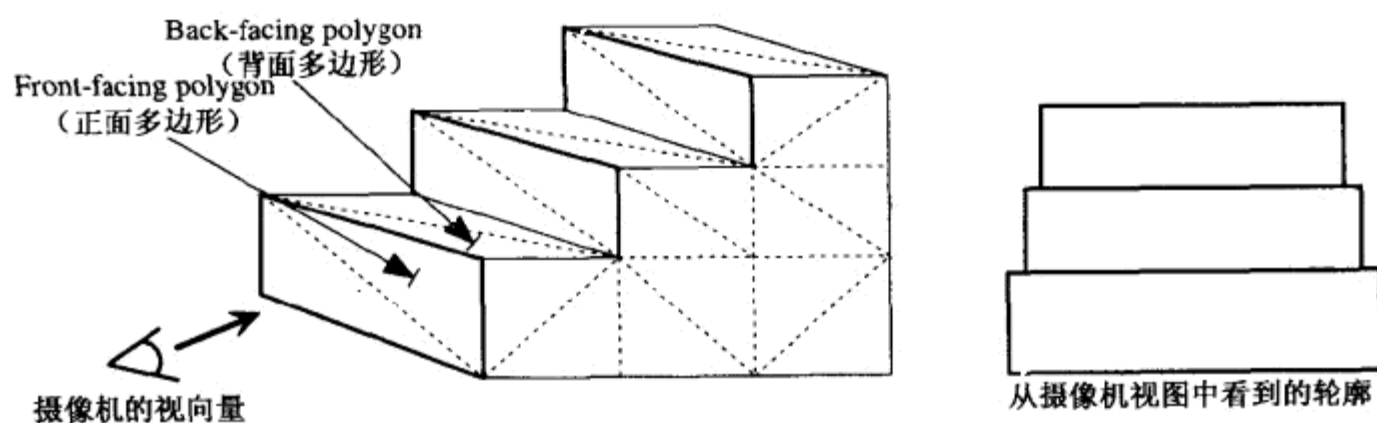


图 5.1.2 左图显示面对第一级台阶检测到的轮廓。右图显示楼梯在摄像机的取景器中显示的轮廓

### 5.1.3 轮廓边缘检测技术

为什么会有这么多不同的轮廓边缘检测技术? 对初学者来说, 考虑不同视频卡上渲染的性能、特点和功用, 几种渲染 API 各有优缺点。本篇描写的技术对许多着墨方法进行了试验, 并加以总结。在不同的硬件和 API 之间, 纹理映射 (texture mapping) 特征、线段渲染以及多边形偏移只是有显著差异的少数几个方面。

### 5.1.4 基于边的着墨

“基于边”一词用来描述分析模型中多边形的边, 检测和渲染重要边的技术。此方法能保证模型完整轮廓以及准确轮廓与折边。其它着墨技术则是通过估计表面曲率来计算轮廓和折边。基于边的技术包括两部分: 预处理 (preprocess) 部分和运行时 (runtime) 部分。

#### 1. 预处理 (preprocess)

预处理的第一步是为模型找到不重复的惟一的边列表, 并为之分配空间。分别检查模型

每个面的 3 条边，如果这些边还没有储存，则把他们插入哈希表储存起来。哈希表建成后，这些边可以存储成一个紧凑的数组。每个边项包括 2 个顶点索引、2 个面索引和 1 个标志项。标志项把边描述为轮廓、折角、边界线和（或）不重要的边。

要精确检测轮廓，需计算面法线。如果模型是静态的（非动态），这些法线可以和折角边一起预计算。对于非渐进网格（non-progress mesh），可以在装载模型或制作模型时处理边界边。对于得到动画物体正确的轮廓，需要权衡的是要在每一帧里重新计算每个面的法线。

## 2. 运行时（runtime）

运行时包括模型重要边的检测和渲染。 $V$  表示观察矢量， $N_1$  和  $N_2$  代表共边三角形的面法线。

运行时着墨（inking）：

（1）如果是动画几何体，则重计算面法线。

（2）对惟一边列表中的每条边：

- 从边顶点位置减去观察矢量位置来计算  $V$ 。
- 若边有一个相邻面或两个不连续的相邻面，则设置边界标志。
- 若为动画几何体，当相邻面的二面角超过已知阈值时设置折皱标志。
- 用  $(N_1 \cdot V) \times (N_2 \cdot V) \leq 0$  检测轮廓，并设置轮廓标志。

（3）遍历处理边列表，同时渲染已设置边标志的边。建立独立的轮廓、折边和边界线列表可优化边渲染，因为每个边列表可在单个 API 调用中渲染。如果所使用的 API 支持多边形偏移和具有宽度线段的渲染，那么它们都能增强最后渲染的视觉效果。

第二步是轮廓检测算法的核心。通过计算面法线和每条边相邻面的观察矢量的点积（dot product），能非常精确地检测轮廓。使用顶点法线（计算轮廓）可以作为一种优化，不过像许多基于顶点和像素着色器的方法一样，这种技术将丢失一些轮廓。参见程序清单 5.1.1 中的代码示例。

优点：

- 跨平台解决方案；
- 完全精确的轮廓边缘检测；
- 美工可为折角边选择阈值；
- 通过  $z$  缓冲解决边的可见性问题。

缺点：

- 仅在特定 API 中可使用具有宽度线段的渲染；
- 必须为轮廓边缘检测计算面法线；
- 线中需要应用偏移，以避免模型轮廓线。

### 程序清单 5.1.1 基于边的轮廓边缘检测

```
// Detect silhouettes
Edge *pEdge;

// Iterate over all edges
for(unsigned int i=0; i < numEdges; i++){
```

```
pEdge = &m_pEdgeList[i];

//Calculate vector from eye
pEdgeVertexPosition =
    pMesh->GetVertexPosition(pEdge->GetVertexIndex(0));

// Subtract eye position from a vertex position on the edge
viewVector = pEyePosition - pEdgeVertexPosition;
viewVector.normalize();

// Calculate the dot products from the face normal and the view vector
uDotProduct1 = DotProduct(viewVector, *pFaceNormalA);
uDotProduct2 = DotProduct(viewVector, *pFaceNormalB);

if((uDotProduct1 * uDotProduct2) <= 0.0f){
    // The edge is a silhouette edge
    AddEdgeToRenderEdges(pEdge, uNumRenderedVertices);
}
}
```

### 5.1.5 可编程顶点着色器着色

可编程顶点着色器对现代图形构架来说是一个极其灵活的图形 API。为什么它对着墨有帮助？我们已经说过，着色依赖于法线和观察矢量的点积（dot product）。可把法线和观察矢量传递给图形硬件，以找到轮廓。

可编程顶点着色器着色运行时（runtime）

- （1）通过图形 API 建立纹理管道。这要求创建或载入一维纹理。装入二维纹理也能工作得很好，不过只使用第一行。贴图过程中，得到的颜色可以从轮廓颜色（通常黑色）过渡到任何想要的颜色。一般建议用一张  $1 \times 32$  或  $1 \times 64$  的 1D 贴图。在贴图坐标  $u=0$  附近的像素颜色逐渐过渡到物体材质的颜色。
- （2）装载可编程顶点着色器的寄存器。需要装载的内容可能依赖于特定的应用以及在何处进行矩阵变换，但通常说来，需要发送世界坐标系到视角坐标系转换矩阵和投影矩阵（projection matrix）。接下来，需要转换顶点位置和法线位置，并且发送给 vertex shader 寄存器，来计算纹理坐标。由于要计算新纹理坐标，不必把这些量传递给图形卡。不过，请记住确保图形单元针对模型生成纹理坐标。
- （3）进行从模型到归一化裁剪空间（homogenous clip space）的顶点位置转换。
- （4）对每个顶点，从视角坐标矢量中减去世界坐标中的顶点位置，来创建观察矢量。计算顶点法线和观察矢量之间的点积，即  $N \cdot V$ ，如程序清单 5.1.2 所示。确保法线和观察矢量位于同一坐标系中，否则结果会不正确。在装载到顶点着色器的常量寄存器之前，可在图形单元或 CPU 中完成。
- （5）存储第 4 步结果作为针对  $u$  纹理坐标的纹理坐标。

另一种轮廓可编程顶点着色技术分两步，第一步沿着其法线挤压顶点。然后，禁止深度

缓冲,并用期望的轮廓颜色渲染模型。接下来,重新激活深度缓冲,并在常规状态渲染模型。

优点:

- 快速,所有的轮廓检测过程都在图形卡上进行;
- 线宽可变。

缺点:

- 欠精确。依赖于顶点法线,而顶点法线又依赖于物体形状;
- 没有检测不含特殊变量的折边或边界线;
- 要求可编程顶点着色器 API 的支持。

#### 程序清单 5.1.2 DirectX8.0 顶点可编程着色的轮廓着色器

```

;-----
; Constants specified by the application
;   c8-c11 = world-view matrix
;   c12-c15 = view matrix
;   c32 = eye vector in world space
;
; Vertex components (as specified in the vertex DECL)
;   v0 = Position
;   v3 = Normal
;-----

;-----
; Vertex transformation
;-----
; Transform to view space (world matrix is identity)
; m4x4 is a 4x4 matrix multiply
m4x4 r9, v0, c8
; Transform to projection space
m4x4 r10, r9, c12
; Store output position
mov oPos, r10

;-----
; Viewing calculation eye dot n
;-----
;first, make a vector from the vertex to the camera
;value in r9 is in world space (see vertex xform above)
sub r2, c32, r9 ;make vector from vertex to camera
; now normalize the vector
; dp3 is a dotproduct operation
mov r3, r2 ;make a temp
dp3 r2.x, r2, r2 ;r2^2
rsq r2.x, r2.x ;1/sqrt(r2^2)
mul r3, r2.x, r3 ;(1/sqrt(r2^2))*r3 = r3
dp3 oT0.x, r3, v3 ;(eye2surface vector) dot surface normal

```

### 5.1.6 高级纹理特征着墨

参考文献[Dietrich00]中详细描述了一种能完全在图形处理器上执行的新技术。这种方法用逐像素  $N \cdot V$  来计算和得到轮廓。用 alpha 测试模式渲染来得到黑色轮廓的颜色，用 alpha 测试的比较值来调整轮廓的线宽。下面列出了基本步骤，要获取更多详细信息，请见参考文献。

#### 基本步骤

- 1) 建立规一化的 (normalize) 立方体贴图。它非常有效地建立包含 6 个正方面的单个纹理。每个面包括一条逐像素 RGB 法线，它代表原点到立方体上这一点的方向。
- 2) 用观察空间法线纹理坐标作为立方体贴图索引。这将输出一个表示  $N$  的 RGB 编码矢量。
- 3) 用观察空间位置纹理坐标作为立方体贴图索引。这将输出一个表示  $V$  的 RGB 编码矢量。
- 4) 将 Alpha 比较模式设为 LESS\_EQUAL。
- 5) 将细线的 Alpha 引用值设为 0，值越大，线越粗。
- 6) 求第 2、3 步得到的颜色值点乘，即  $N \cdot V$ 。<sup>1</sup>
- 7) 对轮廓进行一轮 LESS\_EQUAL alpha 测试，再对着色的物体进行一轮 GREATER alpha 测试。

#### 优点：

- 对适当的硬件，速度非常快；
- 允许线宽变化。

#### 缺点：

- 欠精确。依赖于顶点法线，而顶点法线又依赖于物体形状；
- 需要特定的阈值来检测折角，以及特定的 alpha 比较值来检测边界线；
- 需要特定的硬件来加速。

### 5.1.7 结论

卡通渲染场景中用到的着墨和着色只是非逼真渲染 (nonphotorealistic rendering, NPR) 这个大领域中的两个技术。NPR 中还有许多其它程式化的技术可用来试验，如铅笔素描、水彩、黑白显示、美术渲染，以及技术性显示。着墨是基本工具，非常适合从它开始建立您自己的 NPR 库。使用这个库，可以让您随意改变渲染风格，增强游戏吸引力，从而使您开发的游戏与众不同。

1. 注：把  $N$  点积  $V$  的值作为 alpha 比较值。绘制物体时，做“小于等于”的 alpha 测试渲染，以轮廓颜色来进行第一遍物体的绘制。因为物体轮廓部分  $N$  点积  $V$  的值一般都接近于 0（视角到轮廓顶点的法线方向接近垂直）。所以只有物体轮廓部分被渲染上。第二遍渲染时使用物体正常材质把非轮廓的地方渲染上。

### 5.1.8 参考文献

---

[Dietrich00] Dietrich, Sim. "Cartoon Rendering and Advanced Texture Features of the GeForce 256 Texture Matrix, Projective Textures, Cube Maps, Texture Coordinate Generation and DOTPRODUCT3 Texture Blending." Available online at: [www.nvidia.com/Marketing/Developer/DevRel.nsf/WhitepapersFrame](http://www.nvidia.com/Marketing/Developer/DevRel.nsf/WhitepapersFrame). 2000.

[Gooch01] Gooch, Bruce, and Amy Gooch. Non-Photorealistic Rendering, A.K. Peters, Ltd., 2001.

[Intel00] Intel's Graphics and 3D Technologies Web page. Available online at: <http://developer.intel.com/ial/3dsoftware>. 2000.

[Lake00] Lake, Adam, Carl Marshall, Mark Harris, and Mark Blackstein. "Stylized Rendering Techniques of Real-Time 3D Animation." Non-Photorealistic Animation and Rendering Symposium, pp. 13~20. 2000. <ftp://download.intel.com/ial/3dsoftware/toon.pdf>.

[Markosian97] Markosian, Lee, Michael Kowalski, et al. "Real-Time Nonphotorealistic Rendering. In Proceedings of ACM SIGGRAPH 97, pp. 113~122. 1997.





## 5.2 使用纹理映射的卡通渲染与可编程顶点着色器

Adam Lake, 英特尔架构实验室

adam.t.lake@intel.com

卡通渲染是一种风格性的渲染技术，它可以赋予 3D 场景卡通环境的外观和感觉。本篇文章描述的技术利用了现代实时图形功能，包括纹理映射（texture mapping）和可编程顶点着色技术（programmable vertex shading）。基本思想是使用纹理模拟有限颜色的调色板。为了达到目的，我们修改标准的漫反射着色方程，创建高亮和阴影颜色，并使用这些颜色为每个材质（Material）创建一个纹理贴图。此纹理贴图在实时渲染时，用作颜色查找表。另外，这些技术不需要额外的标记信息——这篇文章描述了纹理图以及每个材质的纹理坐标的创建。

### 5.2.1 卡通着色技术

通过模拟用于动画的着色和着墨赛璐珞（cel）过程来实现卡通着色。这篇文章讲解了着色算法。Carl Marshall 撰写的文章《卡通渲染：实时轮廓边缘检测与渲染》讲解了用于动画的赛璐珞着墨技术。着墨（inking）创建了模型的轮廓，它又由渲染器的上色过程来“填充”。通过检测轮廓、折边以及曲面的材质边界来完成。有关着墨技术的更多知识，可以参考 Carl Marshall 撰写的文章。上色器检查物体的材质性质与光照，并计算用于渲染每个网格的纹理图。上色器也计算曲面的纹理坐标。可以单独应用每篇文章描写的精粹技术，获得独特的效果，或者综合使用模拟动画师使用的完整着墨和上色方法。

### 5.2.2 上色

上色（painting）过程可以分为两个阶段：预处理（preprocess）和运行时（runtime）两个阶段。在预处理阶段，为每种材质分配基于光照和材质性质创建的纹理图。在实时渲染时，此贴图的纹理映射状态被激活。贴图坐标基于顶点法线和光源向量之间的夹角生成。对于点光源（positional light），表面到光源的向量需要对每个顶点进行计算。因为物体表面每个顶点到光源的向量都不一样。所以每个顶点到光源的向量要进行一次开销较大的归一化（normalize）。对于平行光源（directional light），在场景中的所有位置，光源均来自于同一方向。因此，方向相同，表面到光源的向量不

需要对每个顶点进行计算。当性能至关重要时，应该使用平行光源。在大部分情况下，在卡通着色过程中，由点光源提供的真实性用处不大。

### 1. 预处理 (preprocess)

对每一种材质，预处理包含三步。第一步，计算照明漫反射颜色。第二步，计算阴影漫反射颜色。最后一步，基于以上颜色为每种材质创建纹理。在这个例子中，仅需要一维纹理，但使用二维纹理同样起作用，且也许更适合于大部分图形硬件和图形 API 的快速路径 (fast path)。

在下面过程中， $C_i$  指照明漫反射色， $C_s$  指阴影漫反射色， $A_g$  指全局环境光 (global ambient light)， $A_m$  和  $A_l$  分别指材质和光的环境色。

预处理卡通着色：

- (1) 计算照明漫反射色。

$$C_i = A_g \times A_m + A_l \times A_m + D_l \times D_m$$

- (2) 计算阴影漫反射色。

$$C_s = A_g \times A_m + A_l \times A_m$$

- (3) 对每种材质，创建并保存一维纹理图，用两个贴图像素 (texel) 包含第 1 步和第 2 步的结果。假定纹理使用纹理坐标  $u$ ，在纹理尾  $u=1$  处保存  $C_i$ ，在  $u=0$  保存  $C_s$ 。

### 2. 运行时 (runtime)

运行时包括建立具有预处理纹理的纹理管道、计算每个顶点与灯光间向量的点积 (dot product)、使用这个值求每个顶点的纹理坐标。 $L$  指光源向量， $n$  指下面步骤中的顶点法线。

运行时卡通着色：

- (1) 建立纹理管线，使预处理中计算的纹理为激活纹理。设置纹理坐标模式为固定生成的纹理坐标。一般而言，我们对这一步中的多纹理不感兴趣，因此设置纹理混合功能为“替换”。
- (2) 对于每个顶点，计算顶点法线与光源方向向量间的点积。我们还希望钳定此值，不允许为负值，因此得到方程  $\max\{L \cdot n, 0\}$ 。参见程序清单 5.2.1 和图 5.2.1。
- (3) 禁止光照。
- (4) 允许纹理映射。
- (5) 渲染模型。

#### 程序清单 5.2.1 计算 Toon 纹理坐标的函数

```
SetTextureCoordinates(Mesh *pMesh, Light *pLight)
{
    //grab the light.
    //If there are multiple light source in the scene return the //greatest contributor
    pLight->GetLightToWorldMatrix(&lightToWorldMat);

    //convert light direction in light space to light direction in //world space
    vector4 lightDirInLightSpace.set(0,0,-1,0);
    vector4 lightDirInWorldSpace;
```

```

lightDirInWorldSpace = lightToWorldMatrix * lightDirInLightSpace;

//convert light direction in world space to light direction in //model space
matrix4x4 *meshToWorldMatrix = pMesh->GetMeshToWorldMatrix();
vector4 lightDirInModelSpace;
lightDirInModelSpace = meshToWorldMatrix->inverse() * lightDirInWorldSpace;
lightDirInModelSpace.normalize();

int iNumVerts = pMesh->GetNumVerts();

//for demonstration, we assume directional lights, point lights //would need a vector
pointed from eye to vertex for each //dotproduct.

for(int iVert=0;iVert<iNumVerts;iVert++)
{
vector3 vert = mesh->GetVert(iVert);
vector3 normal = mesh->GetNormal(iVert);
float fTexU,fTexV;

//only negate if the light points from light to object
fTexU = -normal->dotproduct(lightDirInModelSpace);
//max(L.n,0)
if(fTexU < 0) fTexU = 0.0;
fTexV = 0.0;
mesh->SetTextureCoord(iVert,fTexU,fTexV);
}
}

```

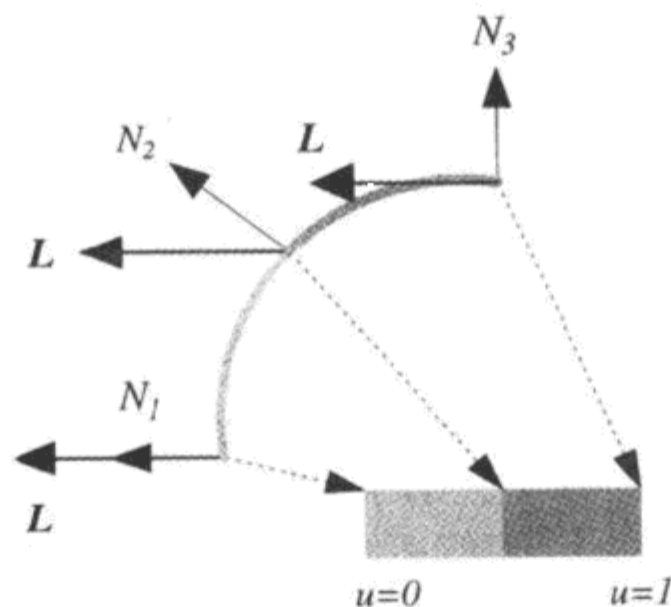


图 5.2.1 在此图中,  $N_1$ 、 $N_2$  和  $N_3$  为表面的法线向量。 $L$  表示平行光源向量。注意, 当  $N_i \cdot L=0$ ,  $i=1$  时, 处于纹理的  $u=0$  处;  $N_i \cdot L=0.5$ ,  $i=2$  时, 为过渡点或纹理中的硬边界。这发生在加亮色与阴影色之间的过渡之时。最后, 若  $N_i \cdot L=1$ ,  $i=3$  时, 到达纹理坐标  $u=1$  处

### 其他方案

前面讲述了一种颜色的一般卡通着色样式，一种是以加亮色，另一种是以阴影色。可以使用许多其他方案来达到各种不同的效果。例如，可以用任意数量的颜色创建纹理图来模拟可用高色板。作为一极端的例子，甚至可以在纹理中使用若干种颜色来近似 Gouraud 着色，我们称之为渐变着色法（gradient shading）。纹理也可以只用黑白质素来填充，这就是黑白着色方式。而且，可以使用权重因子调整  $C_i$  和  $C_s$  来加亮或变暗纹理颜色，获得更多的效果。

使用多纹理，可以增加上述算法的灵活性。通过多纹理技术，在基本卡通着色纹理的基础上，还可以应用其他纹理。例如，通过卡通着色法，可以在创建芳香树脂风格游戏时，让人物造型身穿伪装的裤子。通过多纹理技术，可以建立一个为卡通着色的纹理，另一个针对造型进行伪装的纹理（如图 5.2.2 所示）。



图 5.2.2 在此图中，鸭子模型根据前文所述技术进行着色处理。相关代码参见程序清单 5.2.1

### 5.2.3 可编程顶点着色器

现代图形硬件和图形 API 以惊人的速度在发展。技术的发展使得实现卡通着色器中出现了能有效运用的高性能。可编程顶点着色技术有两个重要优势。一个优势是，允许我们替换光照方程。在传统图形管道，基于顶点的光照方程难以应用到 API 中。这能够针对特定的光照方程进行硬件优化。发展的现代图形架构可以对顶点着色器进行编程。我们可以用任意的照明计算替换传统光照方程。显然，存在依赖于卡的硬件限制，依赖于卡的指令缓存大小、寄存器的数量和类型、指令集等。有望出台规范硬件厂商的标准。优化顶点着色器的软件实现，充分利用 Intel CPU 的 SSE 和 SSE2 指令。

假定已有一个可编程顶点着色器，如何帮助我们进行卡通着色处理呢？如果有硬件，就可以将运行时纹理坐标计算有效地转移到图形单元上。预处理步骤与以前相同。修改运行时块，装载具有法线和光照信息的图形单元上的寄存器，并在卡上执行点积计算。最后，把计

算所得的纹理坐标装载到适当的纹理寄存器。

可编程顶点着色器运行时：

- 1) 建立纹理管线，激活预处理中计算的纹理。因为使用的是可编程顶点着色器，所以不需要关心建立纹理坐标生成模式——在图形单元上完成这一工作。
- 2) 装载可编程顶点着色器的寄存器。需要装载的对象取决于特定的应用以及进行转换的地方，但通常需要发送“世界坐标空间到视角空间”（world to eye）和投影矩阵到寄存器。还需要发送被转换的顶点位置以及顶点法线。同时需要发送用于计算纹理坐标的光源位置。因为在计算新纹理坐标，它们不需要传递到坐标显卡。不过，要记住确保图形单元使用针对该模型的纹理坐标。当处理前面不具备纹理坐标的模型时，容易忽视这个错误。
- 3) 把顶点从世界坐标系空间转换到屏幕空间。
- 4) 对于每个顶点，计算顶点法线和光源方向向量之间的点积  $L_i \cdot N$ 。确保法线和光源向量位于同一个坐标系内，否则得到不正确的结果。完成这一工作的最简单方式是将光源向量从光源坐标空间（Light Space）到模型本地坐标空间（Model Space）。在装载到 Vertex Shader 的常量寄存器中之前，在图形单元中可以完成这一切。
- 5) 保存第 4 步的结果，作为  $u$  方向的贴图坐标。

#### 程序清单 5.2.2 DirectX 8.0 卡通可编程顶点着色器

vs.1.0

```

;-----
; Constants specified by the app-must be loaded using
; dx8Device->SetVertexShaderConstant(reg,&vector,nVectors)
; c0      = ( 0, 0, 0, 0 )
; c8-c11  = world-view matrix
; c12-c15 = view matrix
; c20     = light direction (in model space)
;
; Vertex components (as specified in the vertex DECL)
; v0      = Position
; v3      = Normal
; Instructions
; m4x4    = 4x4 matrix translation, same as 4 dp3s
; mov     = moves values from right to left operator
; dp3     = dot product
;-----
; Vertex transformation
;-----
; Transform to view space (world matrix is identity)
m4x4 r9, v0, c8
; Transform to projection space
m4x4 r10, r9, c12
; Store output position
mov oPos, r10
;-----

```

```

; Lighting calculation - calculate texture coordinate
;-----
; light vector dot vertex normal
; no max? Texturing mode set to CLAMP which clamps for us.
dp3 oT0.x, c20, v3;

```

#### 5.2.4 结论

以上讲解了卡通着色处理的两种方法，还讲解了相应算法中的核心技术。这些技术结合着墨技术，可以用于创建卡通环境，在游戏中得到独特的外观和感受，或获得丰富的经验。使用不同的纹理技术，可以获得其他各种效果，包括数字雕刻、木质雕刻、灰岩、大理石、新闻纸等等。使用本文介绍的基本途径（即纹理映射和可编程顶点着色器），可以实际使用这些效果。随着新硬件和 API 赋予我们新的方式，计算机图形学无疑将在未来的 PC 革命中发挥作用。如果您有兴趣探讨其他渲染方法，可以阅读参考文献。参考文献[Dietrich00]讨论了使用 DirectX 的技术。[Lake00]列出了更详尽的学术性参考文献。[Gooch98]讨论了用于技术图解（technical illustration）方面的技术。[Mark97]讨论用统计法来着墨。愿读者能充分利用这些技术，开发出优秀的游戏！

#### 5.2.5 参考文献

[Dietrich00] Dietrich, Sim, "Cartoon Rendering and Advanced Texture Features of the GeForce256 Texture Matrix, Projective Textures, Cube Maps, Texture Coordinate Generation and DOTPRODUCT3 Texture Blending," available online at [www.nvidia.com/Marketing/Developer/DevRel.nsf/WhitePapersFrame?OpenPage](http://www.nvidia.com/Marketing/Developer/DevRel.nsf/WhitePapersFrame?OpenPage), 2000.

[Gooch98] Gooch, Amy, Bruce Gooch, Peter Shirley, and Elizabeth Cohen. "A Non-Photorealistic Lighting Model for Automatic Technical Illustration." In *Proceedings of ACM SIGGRAPH 98*, pp. 447~452, 1998.

[Gooch01] Gooch, Bruce, and Amy Gooch. *Non-Photorealistic Rendering*, A.K. Peters, Ltd., 2001.

[Intel00] [www.intel.com/ial/3dsoftware/doc.htm](http://www.intel.com/ial/3dsoftware/doc.htm). Contains presentations, papers, and other references to Cartoon Shading, Multi-Resolution Meshes, Subdivision, and Character Animation. There is also an application that demonstrates cartoon shading, pencil sketching, and other NPR techniques, all running on animated multiresolution meshes.

[Lake00] Lake, Adam, Carl Marshall, Mark Harris, Marc Blackstein, "Stylized Rendering Techniques for Scalable Real-Time 3D ACM Animation," In *Symposium of Non-Photorealistic Animation and Rendering (NPAR) 2000*, pp. 13~20, 2000.

[Mark97] Markosian, Lee, Michael Kowalski, Samuel Trychi, Lubomir Bourdev, Daniel Goldstein, and John Hughes. Real-Time Nonphotorealistic Rendering. In *Proceedings of ACM SIGGRAPH 97*, pp. 113~122, 1997.



## 5.3 动态逐像素光照技术

Dan Ginsburg, Dave Gosselin, ATI Research  
ginsburg@alum.wpi.edu, gosselin@ati.com

在游戏中实现动态光照效果的一种常见方法是使用基于顶点的光照技术。基于顶点的光照技术有一个重大的缺陷，为了使逐像素光照效果相近，必须让几何体高度网格化。这篇文章讲解的一些技术可以用于以逐像素为基础执行动态光照效果。这些方法的优点是，它们不要求高度网格化几何体，而且常常能够以少量性能代价在多纹理图形硬件上执行。

### 5.3.1 动态光照贴图的 3D 纹理

下面讲解的第一个动态光照方法是使用 3D 纹理进行动态光照贴图的。通过 3D 纹理与 2D 纹理的关系，可以充分了解 3D 纹理。在传统的 2D 纹理中，在每个顶点中存在两个纹理坐标，它们插值于表面。这两个纹理坐标( $s, t$ )都表示沿着其中一个纹理坐标轴（宽和度）的距离。3D 纹理在 2D 纹理的基础上，通过添加第三个表示纹理深度的纹理坐标( $r$ )来加以扩展。3D 纹理可以看作 2D 纹理片的深度（如图 5.3.1 所示）。纹理坐标( $r$ )用于选择访问哪一个 2D 图，坐标( $s, t$ )与 2D 纹理中的一般用法相同。

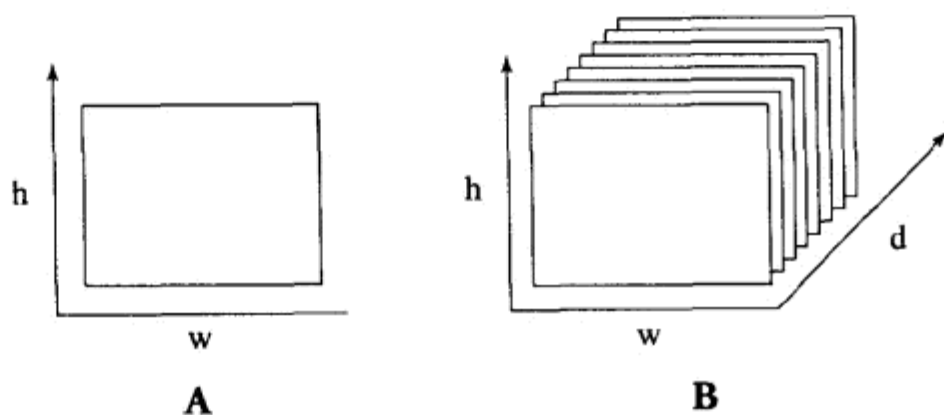


图 5.3.1 A) 2D 纹理; B) 3D 纹理

3D 纹理的一个自然用途是动态光照。游戏中动态光照的一个典型应用是在场景中配备一个照明几何体的移动光源。例如，玩家发射的火箭炮将动态照明静态游戏世界几何体。在 3D 纹理中可以创建任何形式的光源。一个简单的例子是光从中心开始线性衰减（linear falloff）的球形光源。在本文示例中，使用的是具有二次衰减（quadratic falloff）的球形光源（如图 5.3.2 所示）。

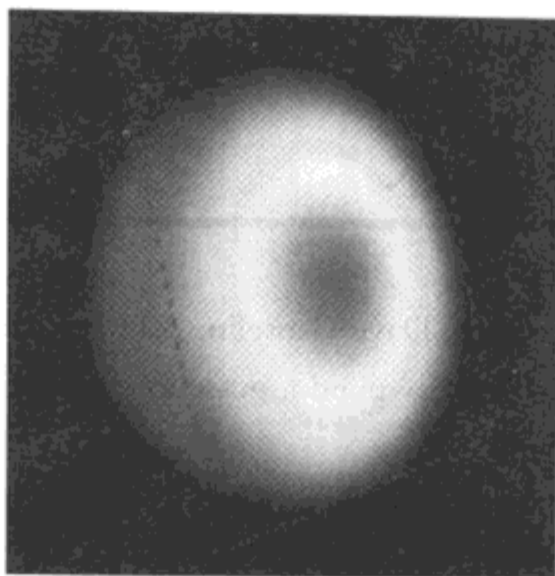


图 5.3.2 3D 光照图的横截面。实际光照图不是单色。这幅图使用颜色编码技术来清晰显示图的衰减。红色亮度最高，蓝色最低

用程序清单 5.3.1 生成 3D 光照图。

程序清单 5.3.1 生成具有二次衰减 3D 光照图的代码

```
for(r = 0; r < MAP_SIZE; r++)
    for(t = 0; t < MAP_SIZE; t++)
        for(s = 0; s < MAP_SIZE; s++)
        {
            float DistSq = s * s + t * t + r * r;

            if(DistSq < RADIUS_SQ)
            {
                float Falloff = (RADIUS_SQ - DistSq) / RADIUS_SQ;

                Falloff *= Falloff;

                LightMap[r * MAP_SIZE * MAP_SIZE + t * MAP_SIZE + s] =
                    255.0f * Falloff;
            }
            else
            {
                LightMap[r * MAP_SIZE * MAP_SIZE + t * MAP_SIZE + s] = 0;
            }
        }
}
```

3D 光照图本身可以在 OpenGL 中使用 EXT\_texture3D 扩展来指定 (OpenGL 1.2.1 部分)。使用 `glTexImage3D()` 把 3D 纹理发送到 OpenGL，它的行为与 `glTexImage2D()` 很像。惟一的区别是指定了图像的深度，而且贴图像素数组包含宽度×高度×深度贴图像素。

例子中的几何体通过基图(b)进行纹理渲染，由 3D 光照图纹理表示的光源(l)进行调整。希望达到的纹理混合为：

$Result = b * l$

使用多重贴图 (ARB-multitexture)，通过 OpenGL 可以配置纹理环境来执行此操作，如下所示：

```
// 3D texture lightmap on stage 0
glActiveTextureARB(GL_TEXTURE0_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_REPLACE);

// Basemap to be modulated by 3D texture
glActiveTextureARB(GL_TEXTURE1_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_MODULATE);
```

已知光位置( $l$ )以及半径( $lr$ ), 3D 光照图中每个顶点( $v$ )的纹理坐标计算如下:

```
s = (v.x - l.x) / lr
t = (v.y - l.y) / lr
r = (v.z - l.z) / lr
```

所幸的是，以上纹理坐标方程在 OpenGL 中使用纹理矩阵很容易计算。首先，定义每个顶点的纹理坐标与其位置相同 ( $s = v.x, t = v.y, r = v.z$ )。接下来，减去光源位置，将每个纹理坐标放入光源空间。给贴图矩阵增加一个减去光源位置的矩阵变换，就可以完成这个操作。现在需要把纹理坐标除以光源半径。在  $s, t$  和  $r$  中缩放纹理矩阵光源半径的倒数倍 ( $1/r$ ) 后，可用纹理矩阵来完成此运算。如下代码块使用统一的比例因子设置上述纹理矩阵：

```
glMatrixMode(GL_TEXTURE);
glLoadIdentity();
glScalef(1 / lr, 1 / lr, 1 / lr);
glTranslatef(-1.x, -1.y, -1.z);
```

请注意，使用不统一的半径，可以修改光源形状，生成非球形光源，而不必修改纹理。例如，让  $x$  上的比例因子大于  $y$  和  $z$  上的比例因子，得到的光源为椭圆体。

经过以上设置,每个多边形用 3D 光照图调整过的基图进行渲染。图 5.3.3 为应用这种 3D 光照贴图技术获得的效果。

注意，虽然 3D 纹理是执行动态每像素点光源（positional light）的自然方法，但 3D 纹理不是内存使用效率最高的技术。在《游戏编程精粹 1》一书中，Dietrich 撰写的文章《衰减贴图》介绍了只使用 2D 和 1D 纹理实现每像素点光源的方法。

### 5.3.2 Dot3 凹凸贴图 ( Bump Mapping )

前面讲述的使用 3D 纹理作为动态光照贴图的效果通过在表面添加逐像素光照扰动(light perturbation)可以得到改进。一些现有图形卡通过 Dot3 凹凸贴图提供了这种功能。在 OpenGL 中, Dot3 凹凸贴图通过 EXT\_texture\_env\_dot3 扩展来实现。这个扩展只是给 EXT\_texture\_env\_combine (GL\_MODULATE, GL\_INTERPOLATE 等)支持的混合模式添加了一个新混合方程。GL\_DOT3\_RGB 方程执行如下计算:

$$\text{Dest.r} = \text{Dest.g} = \text{Dest.b} = 4 * ((A0.r - 0.5) * (A1.r - 0.5) + (A0.g - 0.5) * (A1.g - 0.5) +$$

$$(A0.b - 0.5) * (A1.b - 0.5))$$



图 5.3.3 由 3D 光照图调制基本图

这个混合方程计算任意两种色源的(r, g, b)向量的点积 (dot product)。

为了让颜色值表示[- 1, 1]范围内的向量分量，对每个分量减去 0.5，使颜色值带有符号。不过，执行此减法并且乘以每个分量，得到的颜色值没有正确缩放。为了使点积结果缩放一致，将它乘以 4。

使用这个方程，可以把每个贴图像素的法向量 (N) 编码成纹理的(r, g, b)值。纹理空间中的光源向量(L)可以编码成另一个源的(r, g, b)值。在混合时执行的点积结果得到每个质素的方程  $N * L$ 。这是逐像素凹凸贴图所需要的基本计算。

用于凹凸贴图的纹理应该包含多边形表面上扰动光照的法线。为了生成凹凸贴图，使用一张表示高度场的灰度图。这张灰度图表示原始贴图表面上的凹凸起伏。从这个高度场中，计算每个方向每个质素高度上的变化，就可以生成凹凸贴图。

使用 Sobel 滤镜可以计算高差。Sobel 滤镜提供了检测垂直和水平变化的两种图像滤镜。

$$dx = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix} \quad dy = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}$$

于是，每个贴图像素的法线等于(- dx, - dy, 1)。规一化这个向量，并放到凹凸贴图中每个质素的颜色分量上。图 5.3.4 显示了这个例子生成的凹凸贴图。

在每个顶点，光源向量将保存在顶点基本颜色中。通过平滑着色，光源向量在多边形表

面被线性插值。插值向量用于每个逐素点积 (dot product) 中。在用于点积前, 这个向量必须转换到每个多边形的凹凸贴图纹理空间中。对于点光源 (positional light), 将每个顶点位置减去 3D 光源的中心位置得到每个顶点的自然空间光源向量。在这里, 只要 3D 光照贴图或几何体移动, 就必须重新计算正切空间光源向量。注意, 要使用具有静态平行光光源的凹凸贴图, 正切空间光源向量可以预计算, 因为每个顶点光源向量不会改变。

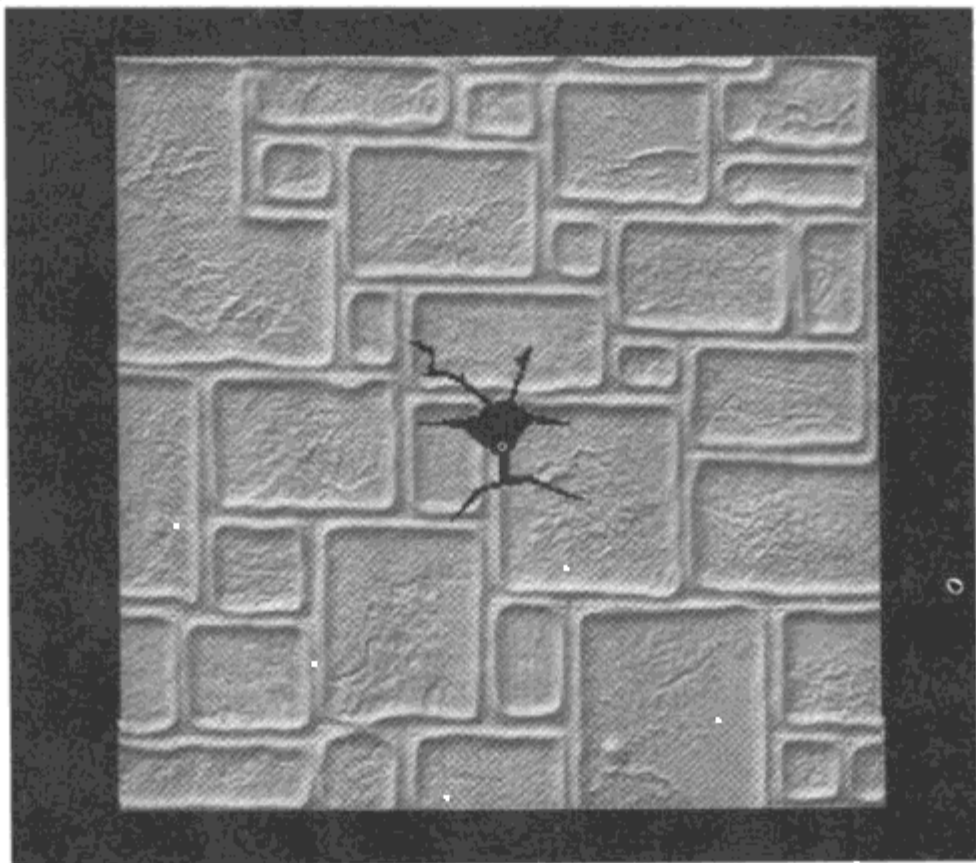


图 5.3.4 使用 Sobel 滤镜从高程场生成的 Dot3 凹凸贴图

为了把光源向量旋转到纹理空间 (即正切空间), 要在每个顶点构建一个局部坐标系。三个向量定义该坐标系: 顶点法线 ( $N$ )、表面切线 ( $T$ ) 以及次法线 ( $B$ )。取包含该顶点的面的所有法线的平均值求得顶点的法线。沿着物体坐标系中的面, 求递增  $s$  或  $t$  纹理坐标的方向, 可以计算得到切向量。取顶点法线和切向量的叉积可以计算得到次法线。在此坐标的基础上, 光源向量可以使用如下旋转矩阵旋转到凹凸贴图空间:

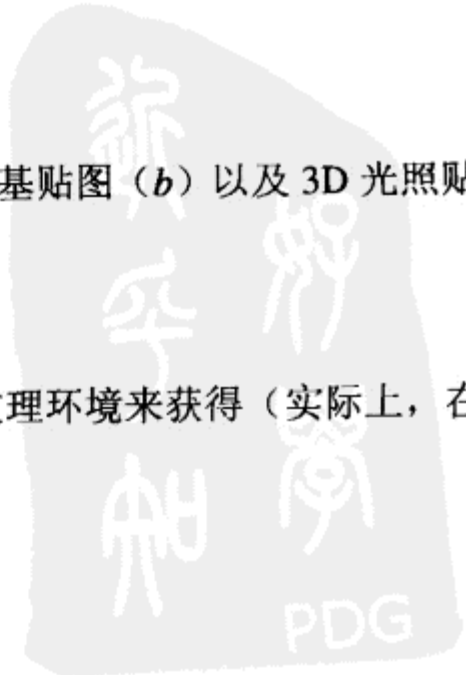
$$\begin{bmatrix} T.x & T.y & T.z & 0 \\ B.x & B.y & B.z & 0 \\ N.x & N.y & N.z & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

最后, 正切空间光源向量保存在顶点的基本颜色中。

已知 Dot3 纹理贴图 ( $dotmap$ )、正切空间光源向量 ( $tl$ )、基贴图 ( $b$ ) 以及 3D 光照贴图 ( $l$ ), 我们希望执行如下纹理混合:

$$Result = (dotmap \text{ DOT3 } tl) * b * l$$

假定硬件有 3 个正交纹理单位, 这可以按如下方式配置纹理环境来获得 (实际上, 在许



多图形卡上, 这一过程必须进行多遍):

```
// dotmap DOT3 t1 on stage 0
glActiveTextureARB(GL_TEXTURE0_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_COMBINE_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_COMBINE_RGB_EXT, GL_DOT3_RGB_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE0_RGB_EXT,
GL_PRIMARY_COLOR_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND0_RGB_EXT, GL_SRC_COLOR);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE1_RGB_EXT, GL_TEXTURE);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND1_RGB_EXT, GL_SRC_COLOR);

// previous * basemap on stage 1
glActiveTextureARB(GL_TEXTURE1_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_COMBINE_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_COMBINE_RGB_EXT, GL_MODULATE);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE0_RGB_EXT, GL_PREVIOUS_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND0_RGB_EXT, GL_SRC_COLOR);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE1_RGB_EXT, GL_TEXTURE);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND1_RGB_EXT, GL_SRC_COLOR);

// previous * lightmap on stage 2
glActiveTextureARB(GL_TEXTURE2_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_COMBINE_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_COMBINE_RGB_EXT, GL_MODULATE);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE0_RGB_EXT, GL_PREVIOUS_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND0_RGB_EXT, GL_SRC_COLOR);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE1_RGB_EXT, GL_TEXTURE);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND1_RGB_EXT, GL_SRC_COLOR);
```

图 5.3.5 显示了应用这一技术的结果。



图 5.3.5 用具有 Dot3 凹凸贴图的 3D 光照贴图调整基贴图



### 5.3.3 使用立方贴图规一化

现在我们已经得到了一个动态照明的逐像素凹凸贴图物体。不过，这个方法有一个主要的缺点。在 Dot3 混合阶段所使用的每个顶点的正切空间光源向量保存在基本颜色中。通过平滑着色把光源向量保存在颜色分量中，该光源向量在多边形表面上被线性插值。但线性插值向量经常不能保持其模长度为 1。图 5.3.6 显示了两个线性插值的 2D 向量。V1 和 V2 长度都为 1，但线性插值后，得到的向量长度小于 1。由于不能保持规一化，所以导致视觉效果下降，尤其当几何体没有被高度网格化时。

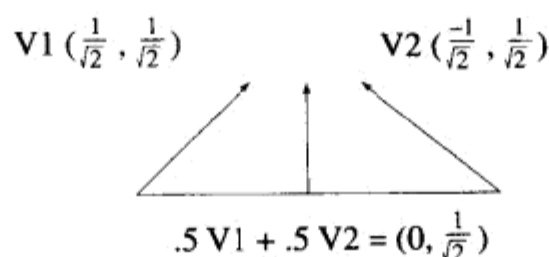


图 5.3.6 两个标准化向量线性插值得到长度小于 1 的向量

为了纠正这个错误，使用立方体贴图来保持光源向量的规一化。每个立方体贴图像素 (RGB) 代表的向量应该与所查找此像素贴图坐标的规一化数值相同<sup>1</sup>。不是将正切空间光源向量保存在主色分量中，而是保存为用于转换到立方贴图的纹理坐标。因为贴图上的每个贴图像素包含纹理的标准化向量，所以光源向量在线性插值时保持标准化。

使用立方贴图的 Dot3 运算可以通过配置纹理环境用两个纹理来完成，如下所示：

```
// Cubemap normalizer on stage 0, just pass it through
glActiveTextureARB(GL_TEXTURE0_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_COMBINE_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_COMBINE_RGB_EXT, GL_REPLACE);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE0_RGB_EXT, GL_TEXTURE);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND0_RGB_EXT, GL_SRC_COLOR);

// dotmap DOT3 cubemap on stage 1
glActiveTextureARB(GL_TEXTURE1_ARB);
glTexEnvf(GL_TEXTURE_ENV, GL_TEXTURE_ENV_MODE, GL_COMBINE_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_COMBINE_RGB_EXT, GL_DOT3_RGB_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE0_RGB_EXT, GL_PREVIOUS_EXT);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND0_RGB_EXT, GL_SRC_COLOR);
glTexEnvf(GL_TEXTURE_ENV, GL_SOURCE1_RGB_EXT, GL_TEXTURE);
glTexEnvf(GL_TEXTURE_ENV, GL_OPERAND1_RGB_EXT, GL_SRC_COLOR);
```

<sup>1</sup> 注：使用要规一化的向量作为 (s,t,r)，来查找立方体贴图像素，得到像素的数值 (r,g,b) 就应该为向量 (s,t,r) 的规一化结果。

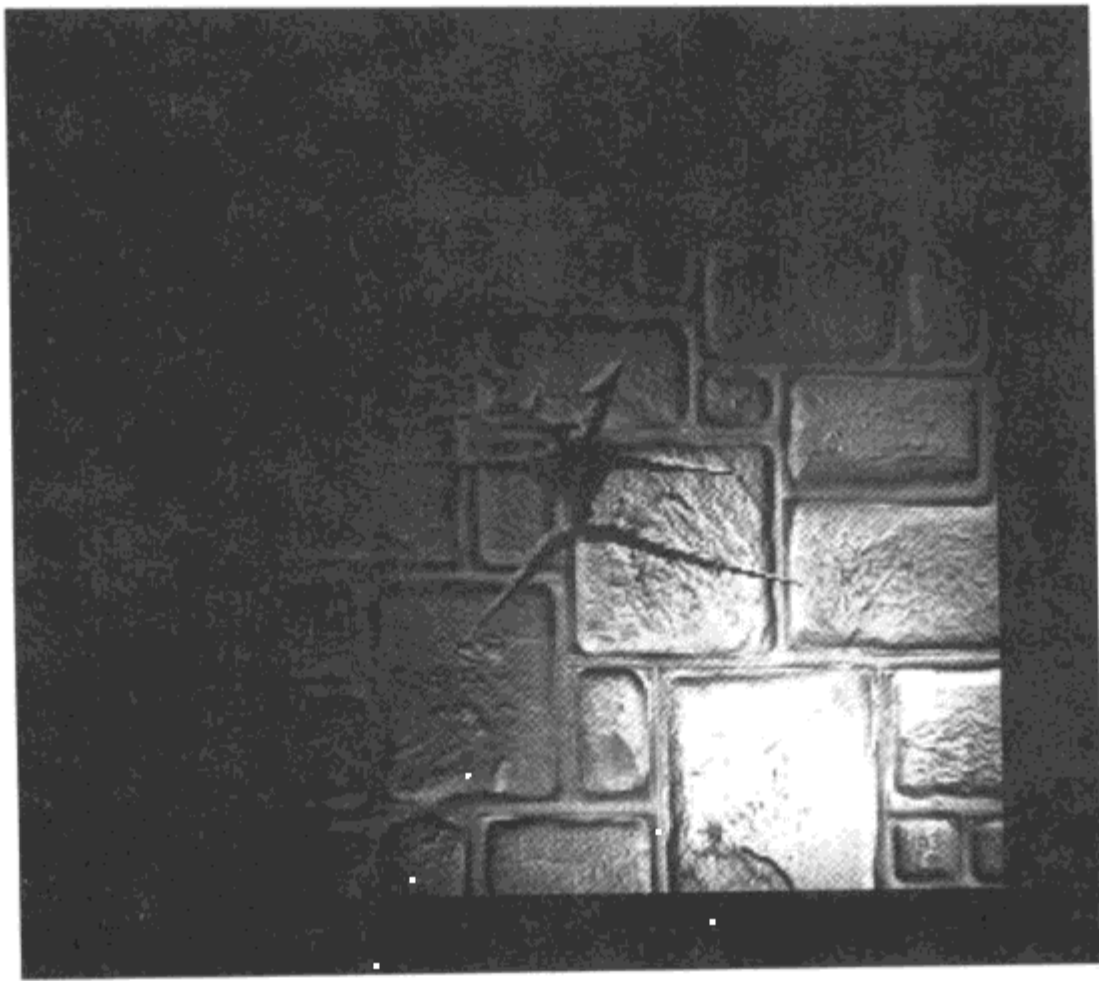


图 5.3.7 用具有 Dot3 凹凸贴图以及立方贴图标准化功能的 3D 光照贴图调整基贴图

#### 5.3.4 逐像素聚光灯 (Per-Pixel Spotlight)

除了用 3D 光照贴图实现的逐像素点光源 (point light) 外, 另一个有用的动态光照技术是逐像素聚光灯 (spotlight)。此时, 不是使用 3D 光照贴图, 而是使用 2D 纹理来表示光源, 并投影到每个面。2D 纹理的内容是包含从中心衰减的聚光灯锥体的横截面。与点光源一样, 自聚光灯的正切空间光源向量放在顶点的颜色分量中。基于锥顶到顶点的衰减因子放在 alpha 通道中。正切空间光照向量将与法线贴图图中的向量一起求 Dot3。插值后顶点颜色的 alpha 值与颜色调制来表现光源随距离的衰减。

Dot3 计算的结果现在需要根据投影聚光灯纹理进行调整。可以配置纹理矩阵来执行投影, 可以把聚光灯投影到面。基于视角位置, 视角方向的和世界垂直方向的向量, 计算观察矩阵, 这样, 首先将纹理矩阵设置为表示从聚光灯角度的视图。观察纹理矩阵最后必须乘以一个表示聚光灯所创建棱锥台的矩阵。于是, 相对聚光灯纹理的纹理坐标就是顶点的坐标系坐标。纹理转换的结果通过 Dot3 计算进行调制, 得到最终光照效果。

#### 5.3.5 参考文献

[Blythe98] Blythe, David, "Advanced Graphics Programming Techniques Using OpenGL" available online at [www.sgi.com/software/opengl/advanced98/notes/notes.html](http://www.sgi.com/software/opengl/advanced98/notes/notes.html)

[DeLoura00] DeLoura, Mark, *Game Programming Gems 2*, Charles River Media, 2000.

Dietrich, Sim, "Attenuation Maps," pp. 543~548.

Dietrich, Sim, "Hardware Bump Mapping," pp. 555~561.

[Mitchell] Mitchell, Jason, "Radeon Bump Mapping Tutorial," available online at [www.ati.com/na/pages/resource\\_centre/dev\\_rel/sdk/RadeonSDK/Html/Tutorials/RadeonBumpMap.html](http://www.ati.com/na/pages/resource_centre/dev_rel/sdk/RadeonSDK/Html/Tutorials/RadeonBumpMap.html)

ARB\_multitexture, EXT\_texture\_env\_combine, and EXT\_texture\_env\_dot3 spec available online at <http://oss.sgi.com/projects/ogl-sample/registry/>



## 5.4 使用 3D 硬件生成过程云彩

Kim Pallister, Intel

kim.pallister@intel.com

我们所玩的大量游戏发生在户外环境中，而且大部分场景都接近于我们所处的真实环境。因此，对于许多游戏开发者，真实渲染地形成一项极具诱惑力，也具有相当难度的任务。但遗憾的是，模拟天空（当然也包括云彩）还没有得到足够的重视，通常只是在事后稍作完善。云彩通常用一层或多层静态卷动的纹理来模拟。初看之下，这样的云彩似乎很好，但只要多看一会儿，它那重复的不自然本质马上便会显露出来。

在这篇文章中，我们将讲解从程序上生成云彩纹理的方法，它能模拟真实云彩的一些性质。另外，因为纹理一般位于图形子系统内存中，因此我们的目标是几乎完全使用图形硬件来生成过程云彩。最后，为了适应更多的机器，我们将介绍提高这种技术的质量和性能要求的一些方法。

### 5.4.1 云彩性质

找出真实云彩的特征后，就可以提供优先特征进行模拟。与所有实时技术一样，可能由于速度的原因需要放弃许多特征，不过以后我们会对其进行处理。

下面是初步观察云彩得到的几点特征：

- 云彩是动画的。它们在天空中随风飘移，不时改变着形状或不断变幻（从云彩的延时拍摄中可以看出）。另外，空中的云彩数量也随时间而变化。这表明需要 3 个变量：模拟速率、风速以及描述其阴暗度的量。
- 云彩形状随时间而改变，小细节变化频繁，而较大的细节需要较长的时间来改变。
- 云彩在结构上显示出极大的自我相似性。
- 云彩覆盖的程度既可能完全覆盖天空，也可能仅有几片孤立的云朵（或无一丝云彩）在清亮的蓝空中。随着云彩覆盖量的变化，云彩的外观也发生变化。在阴暗情况下，云彩变成灰黑色；在清澈天空中，云彩为带有蓝色背景的亮白色。
- 云彩为三维实体。它们被太阳从一边照耀，在另一边投下阴影。在较小比例下，具有云彩自阴影以及各个方向的漫反射光，情况将更复杂。
- 在日出与日落时，云彩的光照也大不一样，因为太阳光是从正交

角度，甚至是从下面照明。

- 云彩趋向于飘浮在一般高度，形成云层。有时存在多云层。因为每个云层是地球曲面之上的恒定高度，所以可以对云层使用某类曲面几何体。

以上是我们从地面上看到的一些特征。在飞行模拟器等应用程序中，我们可以飞越到云彩之上，这将会带来新的难题。本篇文章仍然以地面为视角。

### 5.4.2 生成随机数

几乎与所有程序纹理一样，生成过程云彩的第一步就是生成噪音（noise）。噪音是用于描述自然界中的一种随机要素的术语。噪音要素是已知输入序列（sequence）（如 1, 2, 3...）的函数，生成一个看起来随机的结果序列（如 0.52, -0.13, -0.38...）。这里说“看起来随机”，是因为对已知输入它必须总是产生相同的结果。这使得它成为伪随机，对于已知输入种子值，可以创建相同的结果。另外，我们所讲的杂色通常带有一定的维数（如 2D 杂色、3D 杂色等）。这只是指映射到随机数函数的输入数，与我们查找多维数组一样。其中一个维由某因子（通常是一个素数）缩放，它充分大，可以防止结果的重复性模式显现出显而易见的样式。

生成随机数（或伪随机数）是大量研究的主题。所有方法都在函数复杂性与结果质量之间寻求最佳平衡点。高质量的随机数生成器生成均匀分布的随机数，而且不会长时间重复。

幸运的是，对于本篇文章来讲，使用非常简单的随机数生成器足够了。当后面讲到特定技术时，将对每个像素用不同的输入“种子”值多次调用随机数生成器。累积结果可以掩盖出现的明显重复性。

我们所使用的伪随机数生成器（PRNG）如程序 5.4.1 所示。它使用输入字符  $x$  作为一个多项式中的变量来工作。在这个多项式中，每一项乘以一个素数。通过缩放多项式每一项不同的比例，在它重复前得到一个长序列。符号位被屏蔽掉了，数字被分解到 0~2 的范围，再减去 1，得到 -1~1 的范围。

程序清单 5.4.1 一个简单伪随机数生成器

```
float PRNG( int x)
{
    x = (x<<13)^x;

    int Prime1 = 15731;
    int Prime2 = 789221;
    int Prime3 = 1376312589;

    return (1.0f - ((x * (x*x*Prime1 + Prime2) + Prime3)
    & 7fffffff) / 1073741824.0)
}
```

在使用多个倍频（octave）的情况下，可以构建一个素数数组，并根据倍频程对 Prime1 和 Prime2 取不同的素数，以增加 PRNG 的随机性。对于本篇文章，使用程序清单 5.4.1 中的

素数集足够了。

因为我们要在图形硬件中实现它，而图形硬件不允许执行程序清单 5.4.1 中的所有代码，所以要在  $512 \times 512$  纹理贴图中创建一个函数查找表。在函数重复前，使用  $512 \times 512$  贴图可以得到约 260 000 个数值。我们把这张贴图用作随机数发生器。把这张贴图的一部分拷贝到另外一张目的贴图中，然后使用软件产生随机数来偏移贴图坐标，如图 5.4.1 所示。由硬件把我们的随机数系列（random number series）贴图渲染到目的贴图的四边形中。由此完成对贴图的拷贝。

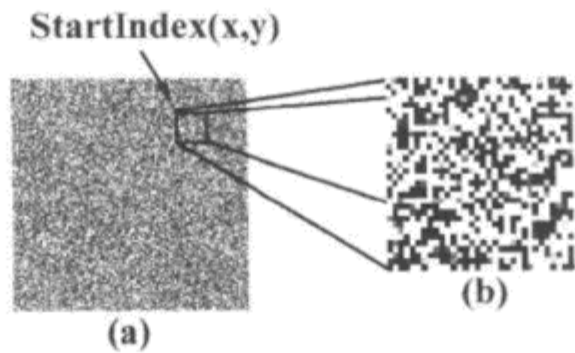


图 5.4.1 (a) 随机数查找表纹理；(b) 生成的 32×32 噪音纹理

受限频带噪音（noise）

Ken Perlin 首先提出了把受限频带噪音（Band-Limited Noise）用作渲染要素的技术。受限频带意味着噪音被限制到一定的频率范围，这可以看作样品间变化的最大变化率。在本文中，这意味着我们希望能够根据频率范围生成对应分布的随机值序列，让中间点在随机样品间插值（interpolation）。只要索引随机数查找表，并用期望的量反采样，使用图形硬件的双线性过滤（filtering）进行插值，就可以完成此工作。去掉高频成份，得到更平滑、更自然过渡的程序要素，如图 5.4.2 所示。

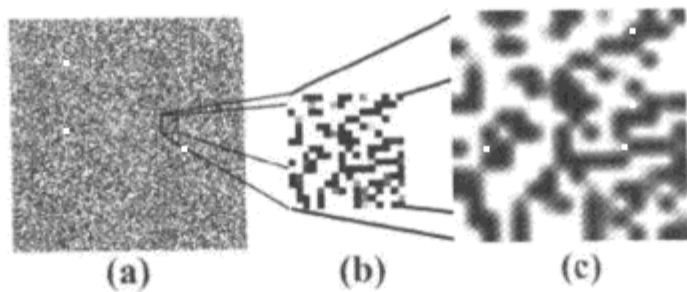


图 5.4.2 (a)采样随机数查找表纹理，创建噪音数组；(b)使用反采样和滤波，创建目标分辨率（即目标频率范围）的受限频带噪音

请注意，与使用噪音创建程序内容不同，使用双线性插值还不够，还必须使用更高级别的滤波，生成更理想的结果。不过，在这里，低质量的双线性滤波可以得到满意的结果。

在这种方式创建的噪音数组得到一个已知频率的单噪音要素。我们将它称为噪音多个倍频（octave of noise），因为以后要将多个倍频（多个相同的频率）结合在一起。不过，首先需要让噪音数组活动起来。



5.4.3 噪音多个倍频的动画

如果希望动画一个噪音倍频程，可以将时间看作第三个维，用这个维来索引随机数生成器。使用图形硬件来实现的方式可以周期性地保留旧噪音纹理、创建新噪音纹理，然后在两者之间进行插值（interpolation），更新到新噪音纹理。更新纹理的速率决定了第三维的频率。

插值过程如图 5.4.3 所示。这是使用线性插值生成人工纹理的一种情形，此时噪音变得更加“集中”于实际采样点。同样，噪音结合到最终效果中后，人工的痕迹不会明显，以后我们会明白这一点。而且，这一点应该记住。如果我们愿意，可以牺牲足够的填充速率，来及时保持多个时间点的倍频，并使用更高级别的插值方案来获得更佳效果。

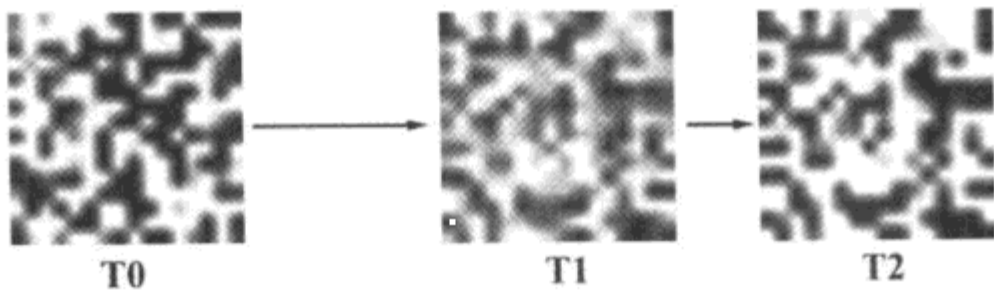


图 5.4.3 插值噪音更新产生噪音倍频程变化。Result = PreviousUpdate \* (1-(T1-T0)/(T2-T0)) + NewUpdate \* (T1-T0)/(T2-T0)。T0 为前一次更新时间，T2 为新更新时间，T1 为两间更新间的当前时间

加合倍频程得到湍流噪音（noise）

生成程序纹理时常用的技术是使用分形和（fractal sum），它是基本噪音函数的比例谐波的总和。这称做分形布朗运动（fractional Brownian motion, fBm）。分形布朗运动用于各种程序渲染技术中，如分形地形、多种类型的程序纹理等等。分形和如方程 5.4.1 所示。在这个方程中，*a* 决定如何访问频谱（*a* 最常见的值为 2，得到 *f*, 2*f*, 4*f*, ...），*b* 决定序列中每一分量的分布域。

$$Noise(x) = \sum_{k=0}^{N-1} \frac{Noise(a^k x)}{b^k}$$

(5.4.1)

方便的是，*a* 取值为 2 既常用于一般程序纹理中，也适合于硬件实现，可以把噪音倍频程（octave of noise）实现为一系列大小为 2 的幂的纹理贴图序列。

2 也既是 *b* 的常用值，也能使实现更容易。使用简单混合执行多遍渲染可以完成噪音倍频程的合成（Composite），如程序清单 5.4.2 中的伪码所示。

程序清单 5.4.2 噪音倍频程合成

```
//note that ScaleFactor is needed to keep results in the
// 0-1 range, it changes {1 + 1/2 + 1/4 + ...} to {1/2 + 1/4 + ...}

float ScaleFactor = 0.5
float FractalSum = 0.0
for i = 1 to NumOctaves
{
```

```
FractalSum = FractalSum*0.5 + ScaleFactor * NoiseOctave(NumOctaves-i)
}
```

图 5.4.4 显示激活更新不同倍频程，并构建合成湍流噪音纹理的过程。

从“水汽”生成云彩

现在得到了动画的纹理贴图、湍流噪音，还需要几步将它变成云彩。理想情况下，将这些噪音用作一些指数函数的输入，得到明显的“浓缩级别”，云彩在此级别之上可见，在此级别之下则不可见。不过，因为在图形芯片中不能使用这种运算，所以必须使用另一种方法。解决的方法有多种。

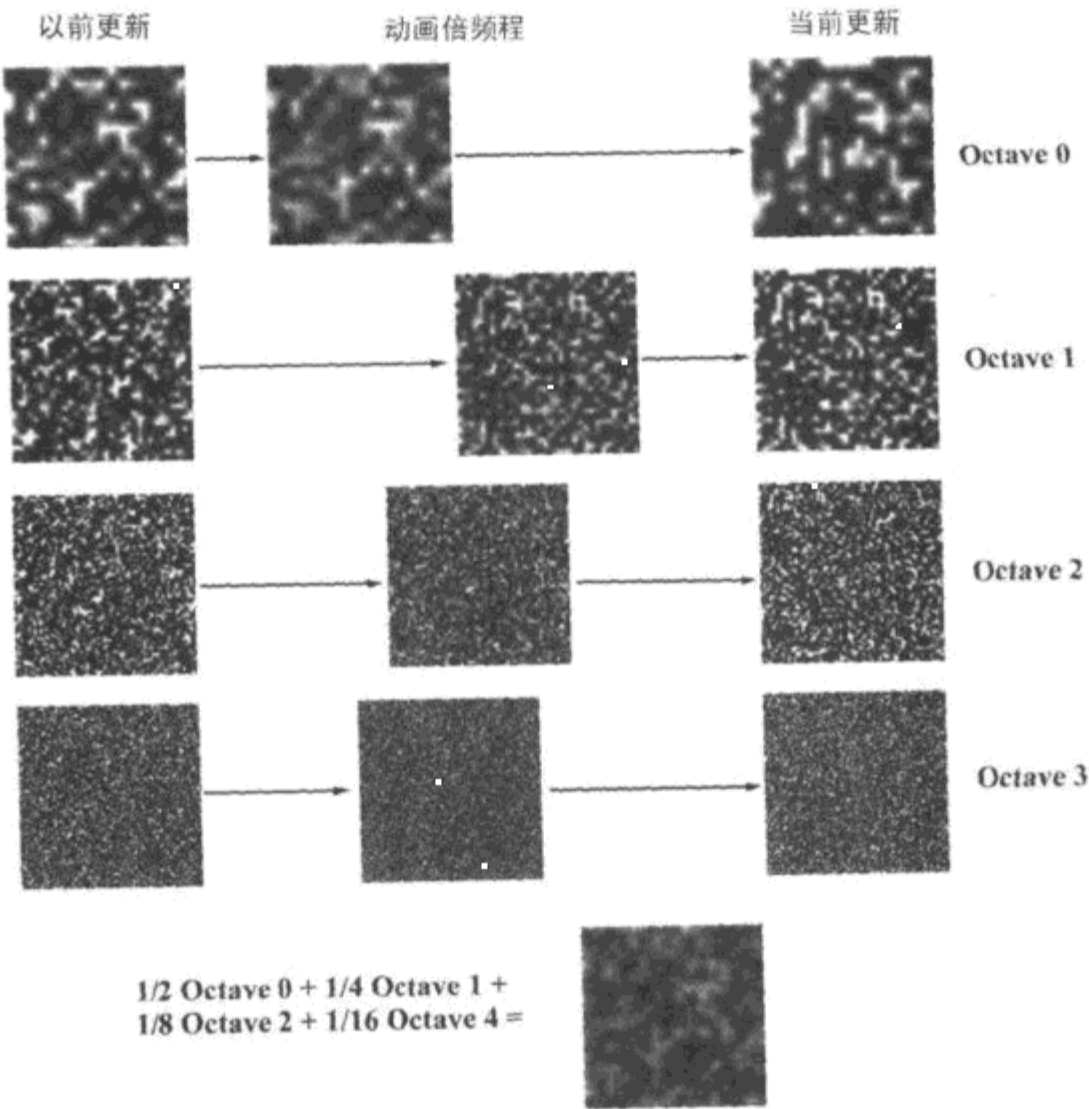


图 5.4.4 动画噪音的几种倍频程的合成

最简单的方式是进行减法运算。把所有贴图像素减去一个固定值。当某像素数值低于某个固定值，则钳位此像素为 0。这样用来使一些连接在一起的云块分离开来，如图 5.4.4 所示。

不幸的是，在保留云彩的同时丢掉了一些灵活性。但可以用顶点颜色和变化值进行过饱和和调制处理，达到弥补的目的。这就是示例代码中所使用的方法。

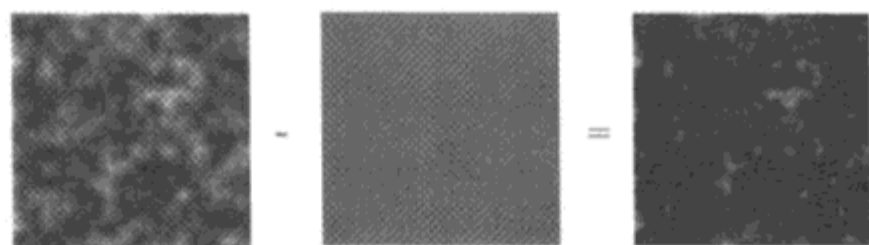


图 5.4.5 减去固定值得到孤立云彩

另一种方式是让到目前为止所提到的所有纹理层具有一个 alpha 通道。这可以像前面的测试中一样进行减法和钳值（但只是在 alpha 通道上），然后执行 alpha 测试屏蔽没有云彩的区域，而不会失去云彩本身中的任何动态区域。不过，这里会发生不同的问题。因为 alpha 测试在任何滤波前完成，所以会发生粗糙的锯齿边缘，除非使用真正高分辨率的目的贴图。

另一种更好的方式仍将使用某类独立的纹理查找，其中，一个纹理的颜色值影响后续阶段获取质素的位置。纹理独立查找的一个例子是 DirectX 下某些硬件支持的 BumpEnvMap 渲染。将来，会有越来越多的硬件支持这类运算。通过这类运算，在纹理贴图中更容易编码指数函数和查找相应的结果。

#### 5.4.4 贴图到天空几何体

图 5.4.5 中已经准备了云彩纹理，下面可以贴图到天空几何体中了。几何体选择取决于您的应用。示例应用使用的是所谓的“skyplane”，它就是矩形的三角形剖分网格，把顶点拽到球半径上。可以想像成把一块布盖在一个大充气球上。它并不十全十美，但使贴图纹理更容易。skyplane 如图 5.4.6 所示。

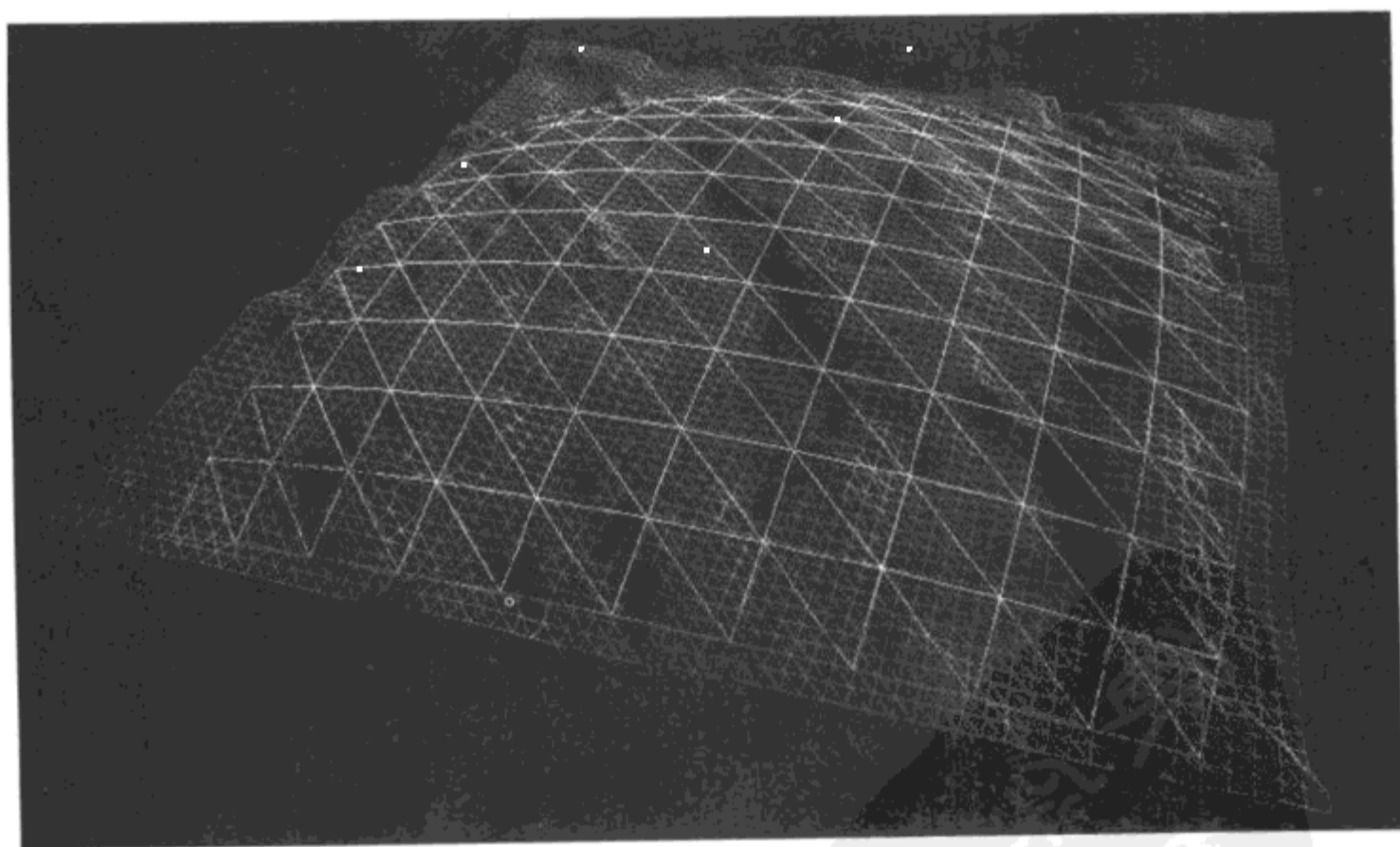


图 5.4.6 地形几何体和用于云彩纹理的 skyplane

### 5.4.5 功能延伸

---

到这一步,根据希望耗用的应用时间和图形硬件填充率,可以对云彩进行更多的处理了。示例代码中完成了一些工作,对于其他列出来的工作,有待读者根据思路自己去实现。可以添加的功能有:

- 噪音驱动风方向、云彩覆盖。使用其他噪音函数修改时间上的其他变量可以得到一些戏剧性的结果,例如,风可以改向,云彩级别可以以小时或天为时间单位增加或减少。
- 让云彩有浮雕效果,赋予它 3D 外观。这需要一遍额外的处理,而且需要对顶点 UV 值进行一些修改。根据太阳的方向,赋值顶点第二套纹理坐标,并进行偏移和缩放。背对太阳一边的云彩变暗,正对太阳一边的云彩变亮。
- 使用最终结果纹理,再加上相减式混合,云彩可以在地面上投下阴影。地形需要另一套纹理坐标,否则为了获得此效果,必须使用纹理投影。
- 修改光照和镜头眩光强度。因为我们知道纹理如何贴图到几何体,也知道太阳的照射角度,所以可以计算对应于视线中太阳位置的准确质素或质素组。可以使用它来更改地形上的光照,或者临时渐变镜头眩光。请注意,当修改光照时,增加云彩覆盖应当减小太阳在场景中的平行光强度,但要增加环境级别,因为真实的云彩将在各个方向漫反射太阳光。

### 5.4.6 硬件限制

---

在这种技术中,让图形硬件承担许多繁琐的工作,与在每一帧认为的锁定和修改贴图相比,可以获得更好的性能。但同时也有一些缺陷。缺陷有三个方面:

- **渲染纹理的支持。**本文描写的技术广泛使用了这一功能。不过,支持纹理表面渲染远未普及。最近发布的 API 版本,更加容易检测是否支持该功能,但是仍然需要一定的技巧。在不支持渲染纹理表面的情况下,可以使用从帧缓冲复制到纹理的办法,但性能有所降低。
- **精度限制。**当前硬件将纹理和渲染目标值保存为整数,通常是每像素 32 bit 或 16 bit。因为我们实质上是在“单色”中工作,生成灰色着色下的纹理,所以被限制到每像素 8 bit,甚至是更坏的情况,即每像素 4 位。后者将产生明显的人工迹象。注意,这意味着高频倍频程只给最终结果贡献了少量的颜色数据。
- **动态范围限制。**因为整数值表示 0~1 范围的数据,所以我们被迫缩放和偏移值来适应这一范围。因为噪音函数最先的返回值在 -1~1 的范围中,所以必须压缩它们以适应可用的范围。这需要做额外的工作,也放大了精度限制产生的人工迹象。
- **硬件指令集和功能的限制。**图形硬件的指令集限制束缚了可以完成功能的范围。如果可以使用幂函数来生成噪音贴图(像指数函数一样),也许将会产生一些有趣的结果。但是我们被限制在简单的算术运算上。另外,我们使用了很多渲染到贴图的技术,这不是在所有硬件上都具备的功能。在将来,随着硬件功能的增强,这个问题会逐渐得到解决。

### 5.4.7 可伸缩性

如果希望在商业游戏项目中使用这个技术，而且面向 PC 这样的目标平台，性能可能在各台机器上不同，因此我们应考虑该技术的可伸缩性。即使是面向固定的平台，游戏的负荷也可能不同。因而，我们可能希望在重负荷情况下，对该技术的质量和性能有一定的保障（具有一定的伸缩性）。

有以下几种方式可伸缩该技术：

- **纹理分辨率。**可以降低不同中间阶段所使用的纹理分辨率，节省内存和填充率消耗。请注意，改变颜色深度意义不大。因为当贴图为 16bit 时，人为加工的痕迹会非常明显。
- **纹理更新率。**不是每张贴图都需要在每一帧更新。尤其时当模拟速度较慢时，就不用在每一帧对贴图进行插值了。
- **使用的倍频程数。**示例代码使用了 4 个噪音倍频程，但使用 3 个噪音倍频程也能得到可接受的结果。相似地，使用 5 个倍频可以在性能较好的机器上取得稍高的质量。

### 5.4.8 结论

图 5.4.7 显示了过程云彩示例程序的最终结果。以云彩作为例子，我们希望这个技术对于读者生成程序纹理有所启示，了解在生成动态程序纹理中如何让图形硬件承担更多的逐像素工作。

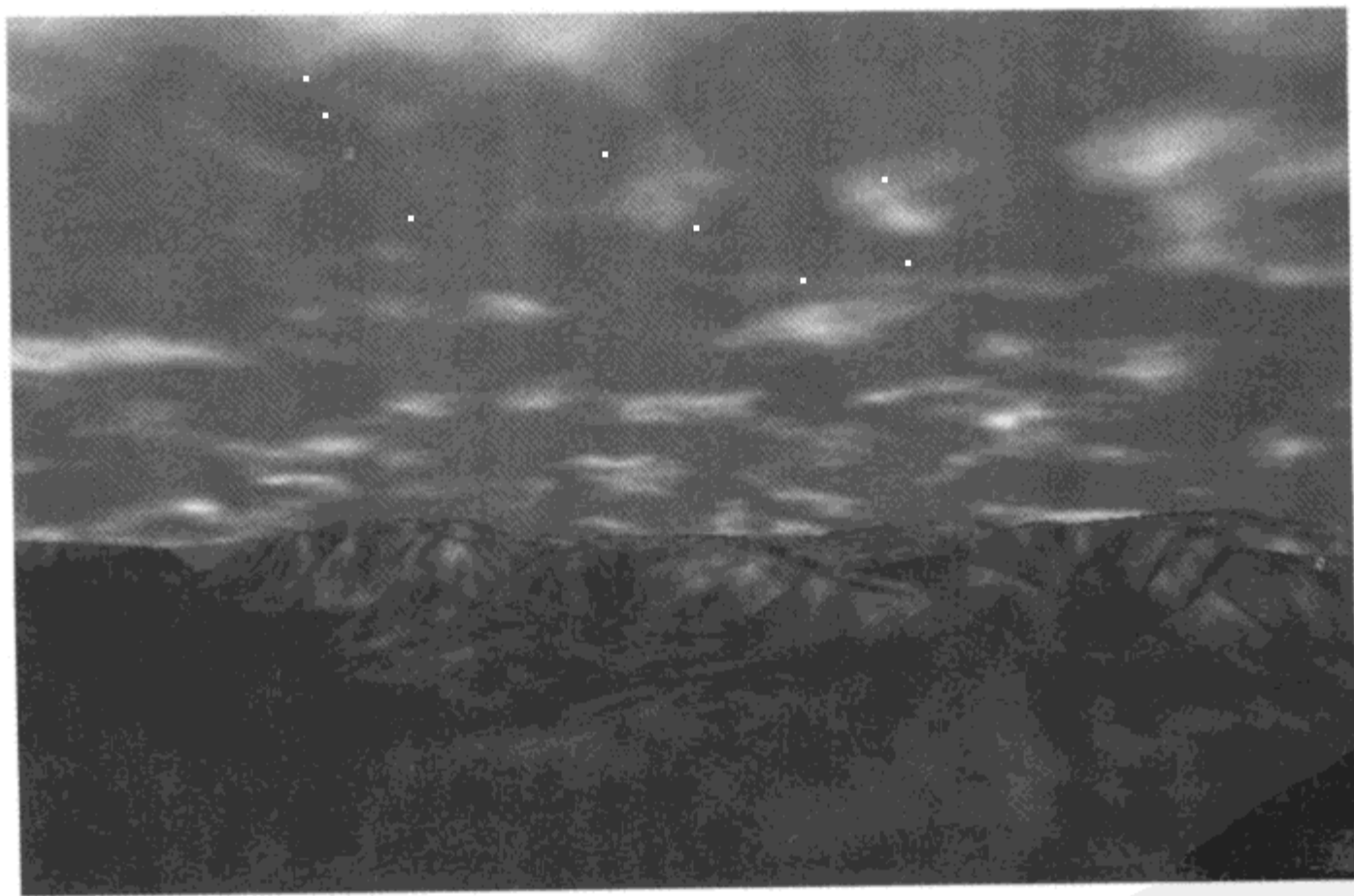


图 5.4.7 过程云彩示例程序的最终结果

随着硬件能力的增强,通过这些技术生成纹理必将成为现实。请读者多多尝试这些技术,并在开发团体中共享成果。

#### 5.4.9 参考文献

---

[Ebert et al, 98] Ebert, David S., et al, *Texturing and Modeling: A Procedural Approach*, AP Professional Inc., 1999. ISBN: 0-12-228730-4.

[Elias99] Elias, Hugo, "Cloud Cover," available online at [http://freespace.virgin.net/hugo.elias/models/m\\_clouds.htm](http://freespace.virgin.net/hugo.elias/models/m_clouds.htm), offers a software-only technique similar to this one.

[Pallister99] Pallister, Kim, "Rendering to Texture Surfaces Using DirectX 7", available online at [www.gamasutra.com/features/19991112/pallister\\_01.htm](http://www.gamasutra.com/features/19991112/pallister_01.htm)





## 5.5 针对较快镜头眩光的纹理屏蔽

---

Chris Maughan, Nvidia Corporation  
cmaughan@nvidia.com

这篇文章介绍一种全新的途径从已经渲染到帧缓冲的像素中生成纹理信息。这种技术可以按几种不同的方式来使用，但这里介绍的方案针对阻塞镜头眩光的常见问题。许多游戏为了准确决定最后场景中的可见物，试图回读帧缓冲中的像素。下面将讲解一种无需 CPU 的协助，也不需要从图形卡读取数据的技术；还将归纳为什么从图形卡回读信息是一种高开销运算，因而应该尽可能避免的原因。

### 5.5.1 镜头眩光遮挡

---

现代游戏为场景添加镜头眩光增加真实感。镜头眩光通常是场景中应用的最后一项，它使用一个渲染为帧的布告版 (Billboard) 的 2D 纹理贴图。进一步完善此技术的办法是，场景中的物体可以阻塞太阳图像。在这种情况下，正确的视觉效果是镜头眩光强度减弱。可以形象地想像成在一条大树成排、阳光灿烂的大道上驱车。如果太阳在树平线之下，随着通过树的视点的改变、到达眼睛的太阳光数量的变化，您就会感觉到有闪烁出现。

检测太阳阻塞的通常途径是首先渲染场景中的物体，包括太阳本身。然后，回读帧缓冲数据，此处正是太阳像素，而且可以推导出场景中的可见太阳光量。完成的方式有两种：可以回读颜色缓冲，并查找匹配太阳颜色的源；或者回读 Z 缓冲，查找同太阳远距离的 Z 值。可见与阻塞太阳像素的比率是镜头眩光强度的近似。现在使用一个 alpha 值设置它的强度，将它混合到最后场景中，就可以绘制出镜头眩光了。

### 5.5.2 硬件问题

---

当在游戏中渲染镜头眩光时通常采用前述方式。虽然它能很好地工作，但在图形管线中导致了不必要的停顿，而这可能会严重降低性能。

现代图形管线非常深。多边形数据不仅在图形芯片管线中排队，而且在大的“分段式”缓冲中排队。一般情况下，游戏中有数千个多边形在分段式缓冲中排队，而且图形芯片在绘制多边形到帧缓冲中时消耗该缓冲。在优秀的并行系统中，游戏可以在 CPU 上进行有用的工作，如物理计算、AI 等，同时图形芯片 (graphic chip, GPU) 排泄分段式缓冲。实际上，如

果要获得系统的最佳性能，推荐大家这样做。图 5.5.1 显示了具有优先级并行性的系统。

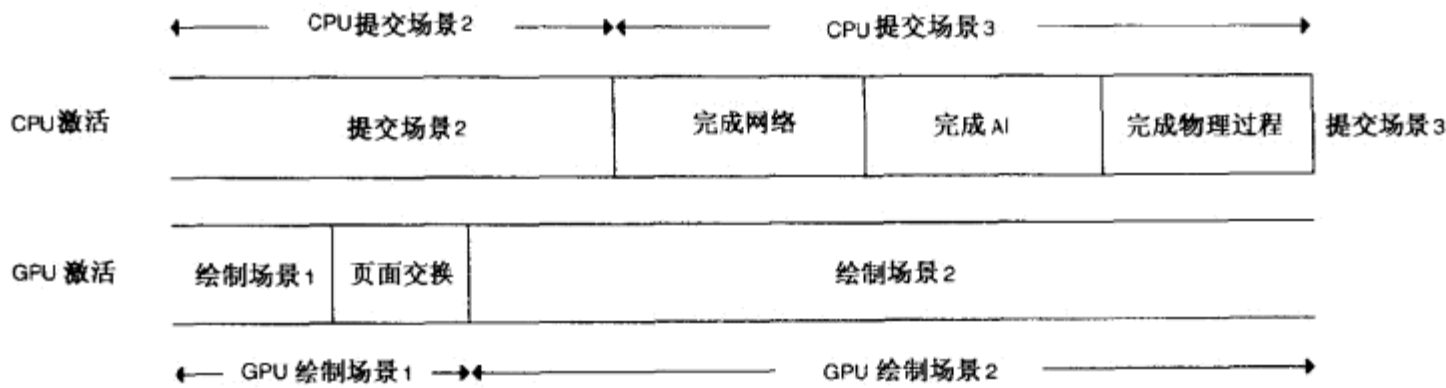


图 5.5.1 理想游戏引擎中 GPU/CPU 的并行性

现在考虑游戏将当前场景中的所有多边形提交给图形芯片后的情形。许多场景仍然在分段式缓冲中排队等候，在一段时间内不会绘制。游戏接下来要做的是回读场景中的内容，决定镜头眩光计算中的太阳可见度。这里就出现第一个问题。为了从帧缓冲读取完整场景，必须等待渲染的完成，而且为了发生此操作，必须刷新整个管线。等待操作发生的同时，CPU 未有效使用，因为它只是在图形卡驱动程序内空闲，等待渲染的完成。理想情况下，我们希望 GPU 和 CPU 时刻同时保持激活。解决这个问题的一种可能方案是，在场景渲染之后，但在计算镜头眩光这前，插入物理/AI 代码。按这种方式，GPU 渲染多边形的同时 CPU 也将保持繁忙。

在这种情况下，如果 CPU 的工作负荷比 GPU 大，就会产生一个问题。因为 GPU 要等待 CPU 完成工作，可以认为系统是非平衡的，它此时可以进行更多的渲染，如图 5.5.2 所示。

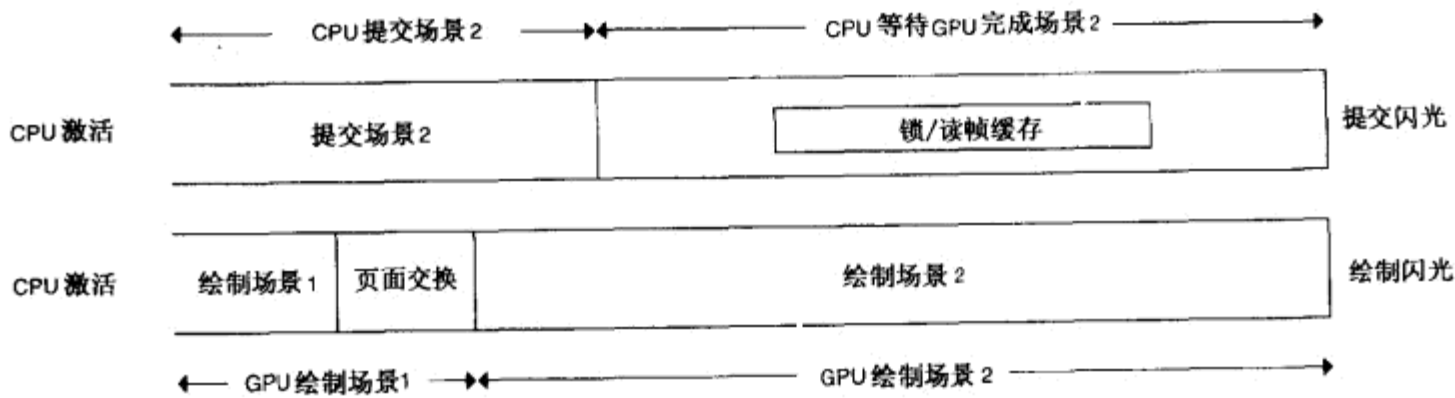


图 5.5.2 渲染结束时的刷新导致的管线延迟

如果 CPU 的工作负荷小于 GPU 的工作负荷，当回读镜头眩光数据时同样会延迟。这不是我们所期望的；我们期望的理想情况是，CPU 在准备游戏下一帧的同时 GPU 完成最后一项工作（如图 5.5.3 所示）。并行是这里的关键所在，插入任何延迟都会造成相当大的损害。

不过，事情并未到此结束。失去并发后，还会出现另外一个问题。从当前生成图形卡回读数据是一个缓慢的运算，在某些情况下根本不能执行。虽然大部分图形卡写入速度为 AGP 2x/4x，但回读速度仍然限制在 PCI 水平上，这是现有 AGP 总线设计的固有本质。结果是，对于 32 bit 帧缓冲，在没有缓冲和突发的情况下，最高只能以 66MHz 的速度读取像素。因此，从图形卡读取 256×256 的区域要花费  $1 / (66\,000\,000 / (256 \times 256)) = \sim 1\text{ ms}$ 。

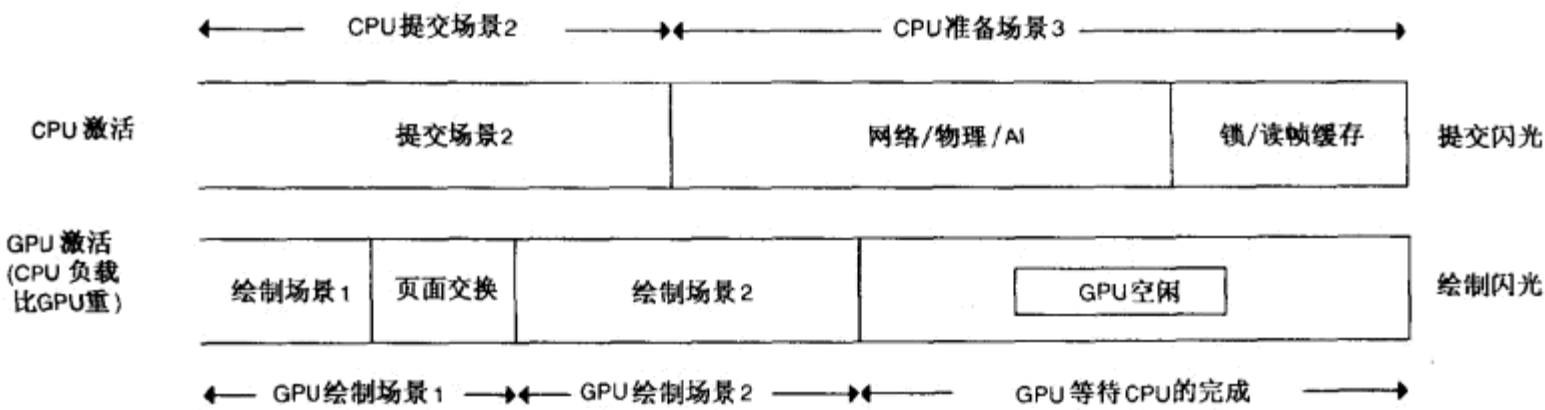


图 5.5.3 CPU 结束工作后，帧结束前的刷新导致的管线延迟

假定 256×256 像素太阳区域投影到屏幕上。如果我们离光源（例如街灯）较近，则投影面积可能变化相当大，直到回读变得非常突出为止。即使我们愿意接受这种性能上的降低，也不能保证图形卡让我们回读数据。事实上，很多图形卡并不允许这样做。大部分图形卡不允许读取 Z 缓冲，有一些图形卡不允许读取帧缓冲，没有图形卡允许读取反锯齿缓冲（antialiased buffer）。

5.5.3 纹理屏蔽

如何才能减轻读取帧缓冲带来的延迟呢？我们通过横向思维来寻找答案。基于前述原因，我们知道，不能从 CPU 向下传递值来调整镜头眩光亮度。现在只剩一种可能性：通过卡中图形内存里已经生成的数据来调整镜头眩光值。这个值必须以纹理贴图中质素的形式出现，而且这个贴图像素需要包含一个照明颜色值，或者包含一个 alpha 值。我们可以在 Direct3D 或 OpenGL 中渲染到纹理。剩下的问题是如何基于一套源像素，实际生成贴图像素。这要通过 GPU 的 alpha 混合功能来实现。alpha 的混合单元基本上是一个加法器，它把 GPU 中的一个值添加到帧缓冲中的某个值上，而且可用于累积颜色信息。再加上场景中的物体用我们选择的颜色进行渲染，就得到合适的解决方案了。下面是达到所需结果的步骤：

第 1 步

创建一个 16×16 贴图像素的纹理贴图，称为“太阳贴图”。我们将在其上渲染太阳，并遮挡太阳的几何体。表面格式应该至少为每颜色分量 8 bit。在一般图形硬件上，这等于 32 bit 的 ARGB 表面。选择太阳贴图为 16×16，因为它包含 256 个贴图像素。我们以后会看到，并不一定要限制到 16×16 大小，如果必要，可以通过滤波技巧来减小大小。

第 2 步

创建 1×1 纹理贴图，称为“亮度贴图”。这是强度数据的最终目标。这只需要单个贴图像素的信息，但当然，这个贴图像素可以按需要变大，以满足特殊的硬件需要。同样，表面的格式应该至少为每颜色分量 8 bit。

### 第 3 步

将包含太阳的场景渲染到太阳贴图。方法有多种，但我们选择一种简单的途径。我们引导观察者直视太阳的中心，同时设置投影矩阵，让太阳填充视口。前端平面正好在眼睛前面，后端平面在太阳之后，用黑色清除太阳贴图。我们在太阳贴图上渲染两种颜色。太阳图像本身用白色渲染。遮挡几何体用黑色渲染，于是覆盖了所有已经渲染的太阳值。如果太阳首先渲染，则不需要 Z 缓冲，因为我们只对阻塞感兴趣，不对深度级感兴趣。太阳贴图现在包含 256 个值，一些值为白色，一些为黑色。在 32 bit 颜色中，这意味着帧缓冲包含 0xFFFFFFFF 或 0x00000000。请注意，为了保持一致，我用设置颜色通道的相同方式来设置 alpha 通道，打开 alpha 混合或颜色混合选项。

### 第 4 步

渲染 256 像素到  $1 \times 1$  强度贴图。如果 API/硬件支持，使用点列表或 point-sprite 来绘制它们，因为它们需要的转换工作比完整的三角形或四边形少。设置纹理坐标，让每个太阳贴图像素被引用一次。设置 alpha 混合单元，在目标亮度贴图上添加每个源太阳贴图的颜色信息。我们还希望从太阳贴图一次只采样一个贴图像素，所以设置纹理采样模式为点采样。注意，我们将使用不同的源贴图像素，多次写入单个目标像素。下一步是混合每个太阳贴图质素与常量值 0x01010101。它实际上是每个颜色通道的  $1/255$ 。结果是，调制后的太阳贴图值要么是 0x01010101（如果源太阳贴图像素为白色），要么是 0（如果源太阳贴图像素为黑色）。我选择白色和黑色作为太阳贴图的值，是因为在渲染中它们一般容易获得，而且可以在演示中显示，便于调试。

这里所使用的技巧是，从太阳贴图向亮度贴图添加 256 个值，在亮度贴图中形成 256 个值的可能范围。亮度贴图上的像素越亮，太阳贴图中就有越多的可见太阳像素；场景本身也是如此。我们遮掩了前面渲染的纹理结果，并将它们加合在一起。我们以后会看到，对不同的应用，纹理屏蔽方式非常有用。

整个过程使用 GPU 从帧缓冲中的信息生成镜头眩光。因为 GPU 管线是串化的，在生成镜头眩光之前，写入了亮度信息，而且管线不需要在前面讨论的运算过程中刷新。另外，因为我们分离了镜头眩光强度与场景渲染，所以可以在任何时候执行这些步骤，而且可以在渲染实际场景期间而不是渲染结束时使用结果值。也可以使用这个值来修改场景的环境强度，来模拟观察者的瞳孔对光的反应。除非希望渲染两次场景，否则不能采用这种需要锁定的方法。

#### 5.5.4 性能考虑

以上介绍的算法看起来有些冗长，但事实上，与读取帧缓冲和刷新管线相比，运算的开销相当少。考虑到在帧尾我们添加了少量的 GPU 工作，它可以继续与引擎的余下部分平行运行。要当心的是，太阳贴图的渲染必须快速。若考虑太阳通常在地面之上，而且对太阳的视场实际上非常小，所以确保这一点。它大量减少了必须考虑的物体，能够快速排除那些通过剖面挑选进行的渲染。如果我们判定某些物体完全阻塞了太阳，则可以进行早期排除，因此

根本不必绘制它们的镜头眩光。另一种选择是，用较少的几何体细节来渲染一些阻塞体。

另一个性能上的问题是，我们使用 256 个混合多边形的渲染来查找强度结果，但这不会有问题，因为现代 GPU 每秒约可以渲染 4 亿个点，消耗的时间只是用于从帧缓冲回读大量数据的时间的一小部分。另外，从帧缓冲回读数据用到 CPU，我们要尽量减少 CPU 在非图形任务中的使用。

### 5.5.5 改进

前述方案能很好地发挥作用，而且得到理想的近似效果。但该技术有一个明显的缺陷：太阳通常是圆形，而我们所使用的太阳贴图是一个正方形纹理。事实上，这不算一个问题，因为我们可以从源太阳贴图按需要采样任何质素，包括圆形采样模式。那当然，为了让它覆盖所要求的采样数，必须让太阳贴图更大。

如果愿意，可以只采样太阳贴图的选择性贴图像素。为此，只要采样太阳贴图中需要的贴图像素，并在写入强度贴图时更改调整运算来放大结果。

注意，我们只是作为一个例子来说明如何解决镜头眩光问题，不过其中的许多技术在其他场合是可重用的。也许强度值可以用于变暗场景来模拟感光过渡效果，或者为人物添加轮廓或晕轮。事实上，只要我们意识到 GPU 具有作为数学解决工具的巨大灵活性，把 GPU 用作一般的平行数学引擎，就可以用多种方式来修改算法。例如，可以使用纹理滤光，一次采样 4 个贴图像素，得到 4 个样品的平均值。在前面的例子中，如果我们依靠双线性滤光得到 4 个样品的平均强度，只需要绘制 64 个点就可得到强度贴图结果。也可以改变混合运算，进行若干数学运算，如使用调制，来缩放值。如果我们想在强度贴图中累积更多的样品，它就可以发挥作用，因为我们能够以动态范围为代价，使用调制给强度贴图结果添加缩放。

### 5.5.6 示例代码



附带光盘上的演示例子说明了锁定技术与纹理屏蔽技术之间的性能差异，也说明了如何实现本文所描述的算法。在菜单中提供了一些选项，可以显示太阳贴图和强度贴图。

当在我的目标系统上运行这个演示时，与纹理屏蔽方法相比，使用锁定方法花费了约 50% 的帧频。帧缓冲变大，性能也得到增加。开放场景的典型数字如下：

纹理屏蔽 (600×500×32bit) = 53.5 fps

帧缓冲锁定 (600×500×32bit) = 27.8fps



源代码中有详细的注释，很容易理解和实战。我使用 Direct3D 8 来编写这个演示，但其概念与 OpenGL 相同。请参考光盘中包含的 README 文件，了解如何运行演示和分析结果的更多信息。README 文件内容还包括如何使用演示中的“Increase CPU Work”选项来研究使用锁定调用与纹理屏蔽技术对系统并行的影响。



### 5.5.7 替代途径

---

为了叙述的完整性，还应提及实现镜头眩光的两种替换途径：

- 有时，可采用一种基于几何体的途径来判断眩光的可见性。这种途径扫描几何体，进行光线交集测试，来决定太阳像素网格的可见性。虽然它能正常工作，但 CPU 的计算时间可能占用较多。而且在使用告示牌时，它不再有用，因为必须扫描源纹理查找“黑洞”。
- 一些图形卡可以异步回读帧缓冲数据。在这种情况下，光源区与其余场景的渲染并行上载。假定在启动回读的阻塞物体被绘制之后、在完成其余场景或光照工作前，场景中存在适当的点，则这一功能可以发挥作用。这个方法当然依赖于 API 的支持，来异步回读数据，而且依赖于硬件的支持来上载数据。在撰写这篇文章时，在 Direct3D 或 OpenGL 中还不提供这种支持，不过，两种 API 的扩展已经得到了提议。

### 5.5.8 参考文献

---

[King00] Yossarian King, “2D Lens Flare,” *Game Programming Gems*, Charles River Media Inc. 2000: pp. 515~518.





## 5.6 实用优先缓冲阴影

D. Sim Dietrich Jr., Nvidia Corporation  
sdietrich@nvidia.com

随着游戏图形在某些方面变得更加复杂，其他方面相比之下显得比较原始。这打破了对于获得虚拟世界中的沉浸式体验（immersive experience）至关重要的一致性。这种现象的一个例子是，详尽逐像素光照的出现，它可以应用到实时整个场景中（参见文章“动态逐像素光照”）。不幸的是，光照的视觉复杂性胜过了许多游戏运用的阴影技术。本篇文章介绍了改进优先缓冲（Priority Buffer）应用的一系列技术，这种优先缓冲首先在文献[Hourcade85]中提出。按照这种方式，我们将探讨其他阴影技术，最后讨论如何创建有用的混合技术。

在实时图形中阴影错综复杂，因为它们是一种场景—图形级别的问题，也许大部分在 CPU 上自然得以解决。不过，为了达到所追求的详尽结果和性能，必须以像素为基础，调节硬件来解决阴影测试。一些图形硬件现在固有支持阴影，包括优先缓冲和阴影深度缓冲。不过，因这些技术不能通过一个单一的标准接口实现，它们都不能使用，所以，我们着重讨论通用技术，通过 Direct3D 或 OpenGL 在现代图形硬件实现。

现代游戏运用的一些阴影技术有：阴影体（Shadow Volume）、深度缓冲阴影以及优先缓冲阴影。每种技术都有不同的优缺点，参见 5.6.1。

阴影深度缓冲[Williams78]是使用 Z 缓冲或表示距离光的逐像素深度的纹理来工作的。我们将着重讲解通过纹理进行工作的技术，最大限度地图形硬件上运用这些技术。阴影深度缓冲工作方式如下：

### 缓冲创建阶段：

对每个灯光

用 0xff 清除深度缓冲贴图

对在灯光世界视锥中的每个物体

计算物体表面到灯光的逐像素深度，深度值从 0x0 到 0xff

使用深度值作为颜色渲染物体到贴图

### 阴影测试阶段：

对每个灯光

创建贴图矩阵，用来把顶点从视角坐标空间移到以灯光为视角的坐标空间

对处在玩家视锥中的每个物体

计算从物体到灯光的逐像素深度

选择深度缓冲为贴图



阴影测试阶段:

对每个灯光

创建贴图矩阵, 用来把顶点从视角坐标空间移到以灯光为视角的坐标空间

对在玩家视野内的每个物体

选择物体 ID 作为固定颜色

选择优先级缓冲作为贴图

对每个顶点

把顶点投影到优先级缓冲贴图上

对物体的每个像素

比较固定 ID 与最近的投影物体的 ID

假如: 固定 ID 大于最近的投影物体的 ID

像素在阴影中

否则:

像素在灯光下

### 5.6.1 比较优先缓冲与深度缓冲

优先缓冲 (priority buffer) 与深度缓冲是非常相似的技术——甚至可以使用相同的逐像素着色程序或纹理阶段设置来执行任何一种方法。

深度缓冲有一个优点: 每个像素被单独对待, 它们均有自己的距离光的深度值。这允许该技术实现自我阴影。

深度缓冲阴影的一个不足之处是, 距离光源的深度一般针对光源的整个影响范围进行计算。如果使用 8 bit 颜色或 alpha 通道来表示前一个例子中距离光源的深度, 则对于整个光源范围只可用 8 bit。对于许多场景, 该精度不足以支持适当的交互对象阴影与物体内的自我阴影 (如图 5.6.1 所示)。



图 5.6.1 距离光源的深度

优先缓冲基于距光源的距离, 在排列顺序的基础上, 给每个物体分配自己的 ID, 达到克服以上困难的目的。

图 5.6.2 是一个以光源为视角, 用不同的 ID 渲染两把椅子的例子。背景通常为白色, 但这里清除成黑色, 更突出显示椅子。

不使用专业的硬件, 有几种办法可克服自我阴影问题。一种办法是将模型分解成凸块, 也许要分解到每一层次的动画模型, 甚至是单个三角形。对于静态几何体的一种自然途径是使用 OctTree 或 BSP 节点将自然世界分解成凸块, 并给每个节点分配一个 ID。不幸的是, 这种途径带来了一个严重的问题, 即锯齿化。

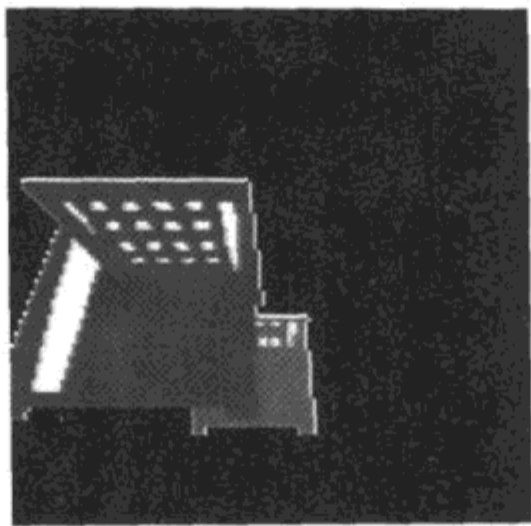


图 5.6.2 光源视角的优先缓冲

### 5.6.2 解决锯齿化问题

当两个具有不同物体 ID 的非重叠物体相互遮光时，产生的锯齿最为明显。之所以发生的简单原因是，后部缓冲和优先缓冲纹理为不同的分辨率和不同的方向。这将导致优先缓冲的一个点样品落在期望取样位置的质素内某处。有时，采样选择期望的贴图像素，有时选择它的邻居（如图 5.6.3 所示）。

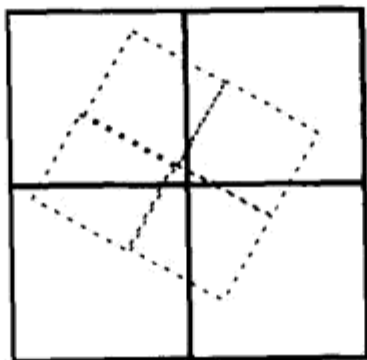


图 5.6.3 虚线显示了与优先缓冲一致的纹理取样位置。大小和位置的不匹配导致了锯齿

通过计算  $2 \times 2$  贴图像素覆盖区的加权平均值，我们将双线性滤波（filtering）用于纹理贴图，以减少这种类型的锯齿。不幸的是，平均物体 ID 无意义。对于物体 ID 100 和 200 的贴图像素，如果平均它们，得到物体 ID 150。这个物体 ID 可能表示一个完全不同的物体，因此，直接的双线性滤波不能完成此任务。

解决锯齿问题的一种途径是偏移纹理一个贴图像素，让样品准确落在  $2 \times 2$  贴图像素区域内。然后，执行阴影测试 4 次，仅当 4 次样品都使像素位于阴影中时才允许此阴影。这解决了锯齿问题，但以 4 次纹理读取以及多遍双纹理硬件处理为代价。

对这种技术的一个改进是“预热”优先缓冲。按前述方法创建优先缓冲后，将 ObjectID 放在优先缓冲的 alpha 通道，执行另一遍渲染到纹理，得到大小与原始优先缓冲相同的纹理。通过这一遍，我们通过取得优先缓冲中与原始样品相邻的 4 个贴图像素，并把它们复制到目标颜色的 R、G、和 A 中，达到“预热”的目的。

按这种方式，在阴影测试阶段，“预热优先缓冲”的单个纹理获取可以提到所有 4 个相

邻样品。通过少量逐像素的数学运算，所有4个阴影测试可同时计算，它们的结果得以综合，锯齿得以消除。如果所有4个测试均一致位于阴影，则只是遮光处理这个像素，如图5.6.4所示。

“预热”技术很有用，因为每一帧只有一步预热，但有许多阴影测试运算，因此它利用了较快渲染的特性。

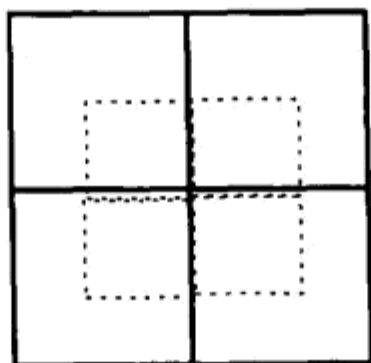


图 5.6.4 虚线显示了与优先缓冲的4个最近邻居一致的纹理样品位置

图5.6.7是图5.6.4的一个“预热”式版本。注意，椅子如何具有有色的轮廓，它表明预热后，R、G和B通道表示不同的ID。

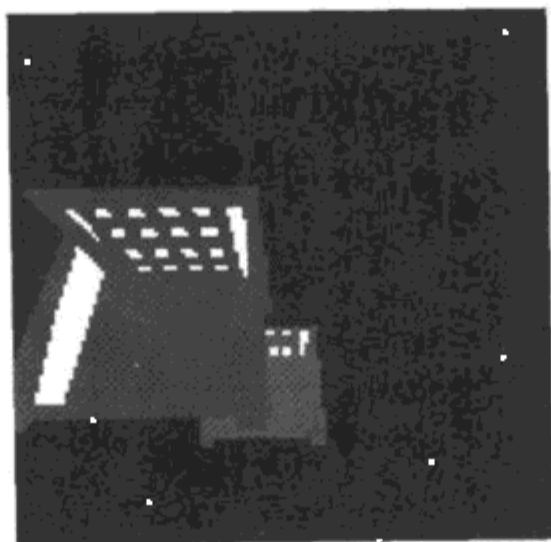


图 5.6.5 光源视角的预热优先缓冲

### 5.6.3 混合途径

也许优先缓冲的最佳用途是用于不相交物体的相互物体阴影，因此不会有锯齿，同时综合使用另一种技术来处理自我阴影。因此，如果两个物体足够远，它们的包围球没有相交，则使用优先缓冲利用较近的一个物体来遮光较远的一个物体。在包围球相交的情况下，有可能相互间产生混淆，因此，它们使用相同的ID，而且必须通过自我阴影测试产生阴影。

结合优先缓冲进行相互物体阴影以及深度缓冲进行自我阴影的一种途径是在纹理垂直方向，将8bit的物体ID编入红色通道中，在纹理水平方向，把物体8bit的深度值编码到绿色通道中。这里采用的思路是，将8bit深度缓冲仅用于单个物体内的自我阴影，因此，8bit可以分布到单个物体的范围。这解决了标准深度缓冲阴影的一个主要问题，即光源范围内的

精度限制问题（如图 5.6.6 所示）。

执行这一双重阴影测试的技术是减法和点积，并复制到 alpha 通道中。alpha 测试可用于判断某像素是否位于阴影中，因而允许这一技术工作于具有双重纹理和点积运算的任何硬件上。

另一种混合途径是仅针对自我阴影使用模版阴影体。在计算了阴影体的边后，可以只把阴影体的边挤压出到物体的包围盒，而不是扩展到光源的范围限制内。这减小了模版阴影体技术的一个潜在性能缺陷——当阴影体穿过摄像机视平面时，像素填充率的耗费。将它们限制为只允许自阴影，也具有如下好处，即完全忽略可视锥（viewing frustum）之外的物体的自我阴影测试。

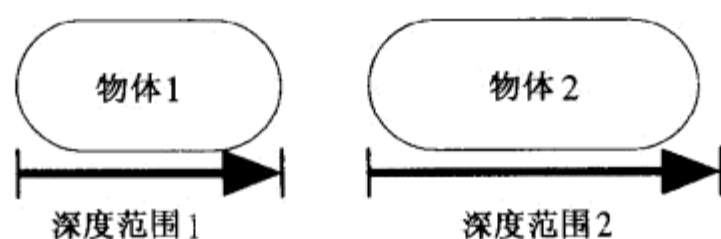


图 5.6.6 相交每个物体的深度量

通过优先缓冲获得自我阴影的另一种途径，是在 CPU 上执行每个顶点到光源的光线投影，同时检查光线与物体的相交。如果光线相交物体的实心部分，则将顶点漫反射色的 alpha 通道设置为 0。在光照过程中，通过光的光照分布调制这个 alpha 值。如前所述，可以只针对可视锥中的物体来执行，甚至对远处的物体可以合并到一起。这种途径的变种是，只在模型的各个块（如肩与臂）上执行测试，也就是说，对每个模型层次级别上执行一次测试。这就为远处和不重要的模型提供了一个理想的细节层次选项。

#### 5.6.4 小结

优先缓冲（priority buffer）是得到相互物体阴影的一种理想途径，它克服了光源范围的精度限制问题。虽然标准优先缓冲对于处理自我阴影不是理想之选，但是可以结合使用其他技术来获得实时游戏中的完全阴影解决方案。

#### 5.6.5 参考文献

[Dietrich00] Dietrich, D. Sim, “GDC 2001 Presentation—Shadow Techniques,” [www.nvidia.com/marketing/developer/devrel.nsf/bookmark/0419FBACDE043ABF88256A1800664C06](http://www.nvidia.com/marketing/developer/devrel.nsf/bookmark/0419FBACDE043ABF88256A1800664C06).

[Hourcade85] Hourcade, J.C., and A. Nicolas, “Algorithms for Antialiased Cast Shadows,” *Computers and Graphics* vol. 9, no. 3, pp. 259~265, 1985.

[Williams78] Williams, Lance, “Casting Curved Shadows on Curved Surfaces,” *Computer Graphics (SIGGRAPH '78 Proceedings)*, vol. 12, no.3, pp. 270~274, August 1978.



## 5.7 替用体技术：添加点缀

Tom Forsyth, Mucky Foot Productions

tomf@muckyfoot.com

恐怕许多读者对替用技术（Impostoring）这一术语并不熟悉。不过，这一概念在 3D 图形历史上却曾以不同形式多次出现。简单地说，它与 3D 场景中使用精灵（sprite）有关，但这种子画面不是一种事先美工绘制或者渲染精灵，而是在运行随时得到更新。事实上，当渲染如用于 VR 穿行的城市风景时，就开始出现替用技术的学术研究了。不过，游戏对动态更新版本更感兴趣，因为根据游戏的本质，它们所处理的场景是不断变化的。

复杂三角形物体并不是每帧都进行渲染，而只是偶而渲染到纹理中，通常每 5~50 帧渲染一次。对每一帧，纹理贴图到简单得多的三角形物体上，用以取代复杂物体的绘制。

替换的主要目标是包含许多小型静态物体的场景，即点缀（clutter）。每个这样的物体都将使用一个替换体，而且大部分都将以一个比主场景帧频低很多的速率进行重绘。这样，感觉到的三角形密度比实际密度高得多，而且所有这些偶然的物体得到的视觉质量大大增加了场景真实性。计算机游戏中的办公室或房子与真实物体之间的主要区别是存在点缀的数量。

虽然新卡拥有庞大的三角形吞吐量，使用替换体仍然有着巨大的优势。总线带宽（在 PC 中就是 AGP 总线）常常是一个瓶颈，而且对某些物体减小总线带宽后，就可以在其他物体上使用更高的三角形细节。一种替换体就是单个纹理，而渲染该物体通常需要多个纹理以及多纹理层——更改纹理刷新硬件的纹理缓存，也可能需要驱动程序（API 或应用程序）的额外纹理管理。绘制每帧的物体要求它照亮每一帧，即使照明没有变化，而且随着照明技术变得更加复杂，照明的开销越来越大。最后，往往存在大量与物体绘制相关的应用开销，甚至在单个 API 调用之前也是如此，使用替换体可以避免许多这类工作。

### 5.7.1 整个过程

替换包括几个关键部分，分别为：

- 每一帧中，在屏幕上渲染 impostor 贴图。我们称这为“渲染”替换体。
- 以低速率渲染替换体的网格到纹理。我们称之为“更新”Impostor。

- 决定何时更新替换体，以及何时不更新。
- 有效地利用硬件。

### 5.7.2 渲染替换体

#### 1. ONE:INVSRCALPHA 样式

一个替换体实质是一幅颜色图像，它具有定义图像透明度的 alpha 通道。已知对于这类混合有两种基本选择：预乘 alpha (premultiplied alpha) 和常规混合 alpha。

常规 alpha 为标准的 SRCALPHA:INVSRCALPHA 样式。其他样式为预乘 alpha，即 ONE:INVSRCALPHA 样式。

使用哪一种样式，取决于渲染 alpha 混合物体到替换体，再渲染替换体到屏幕上时，哪一种样式能产生正确的结果。将像素 P 渲染到一个替换体纹理（它被清除到黑色），产生替换体像素 I，在现有帧缓冲像素 F 上渲染，生成结果像素 R。

如果使用非预乘 alpha，期望结果为：

$$R = P \times P_{\alpha} + F \times (1 - P_{\alpha})$$

渲染到替换体得到结果：

$$I = P \times P_{\alpha} + 0 \times (1 - P_{\alpha})$$

$$I_{\alpha} = P_{\alpha} \times P_{\alpha} + 0 \times (1 - P_{\alpha})$$

渲染到帧缓冲，得到：

$$\begin{aligned} R &= I \times I_{\alpha} + F \times (1 - I_{\alpha}) \\ &= P \times P_{\alpha} \times P_{\alpha} \times P_{\alpha} + F \times (1 - P_{\alpha} \times P_{\alpha}) \end{aligned}$$

它的价值不高，而且与我们期望的结果极不相似。通过预乘 alpha，期望结果为：

$$R = P + F \times (1 - P_{\alpha})$$

渲染到替换体，得到如下结果：

$$I = P + 0 \times (1 - P_{\alpha})$$

$$I_{\alpha} = P_{\alpha} + 0 \times (1 - P_{\alpha})$$

渲染到帧缓冲，得到：

$$\begin{aligned} R &= I + F \times (1 - I_{\alpha}) \\ &= P + F \times (1 - P_{\alpha}) \end{aligned}$$

此结果令人满意。许多应用没有使用预乘 alpha，但要适应现有引擎来使用它相当简单。

要注意的一件事是，现在 alpha 通道结果必须精确，即使渲染不透明三角形时也是如此，因为该结果将被写入替换体，而且将影响替换体渲染到屏幕时的结果。非 alpha 混合（即不透明）渲染到替换体必须确保 alpha 结果为 1，否则，背景可能会透过。幸运的是，完成以上工作非常容易，但相对于常规工作，却需要更加小心。

替换技术的一个问题是，alpha 混合效果渲染到替换体中必须能够在预乘 alpha (premultiplied alpha) 方案内表达其效果。这意味着，乘法颜色混合（主要是 alpha 混合中具

有 COLOR 参数的对象) 不会如期发挥作用, 因为替换体只有一个 alpha 通道来表达允许背景透过的量。所幸的是, 这类效果很少用于适合于替换处理的物体上。注意, 它只适用于实际透明效果。只要能严格控制最终的 alpha 通道值, 使用 alpha 混合综合各遍次 (光照贴图、细节贴图等) 的多遍渲染也不错。

## 2. 广告牌四边形 ( Billboard Quad )

渲染替换体的最明显方式是只渲染表示替换物体的四边形。只要没有物体移动, 也没有摄像机移动, 就能得到很好的效果。

不幸的是, 像素并不是应用在 3D 场景中必须考虑的全部, 还必须考虑深度。四边形只能表示单深度平面, 但它表示的物体将覆盖多重深度。

因此, 对于四边形, 应用需要决定绘制的深度。但对于一些最常见的应用, 没有合适的单深度。如图 5.7.1 所示, 替换体位于两面墙的中间。它的脚伸入邻近的墙体, 它的头消失在远处的墙中。如果替换体与其他物体靠近, 则四边形不理想。对于飞行的物体, 通常不会太接近物体 (例如在飞行模拟中), 它们可能非常理想。

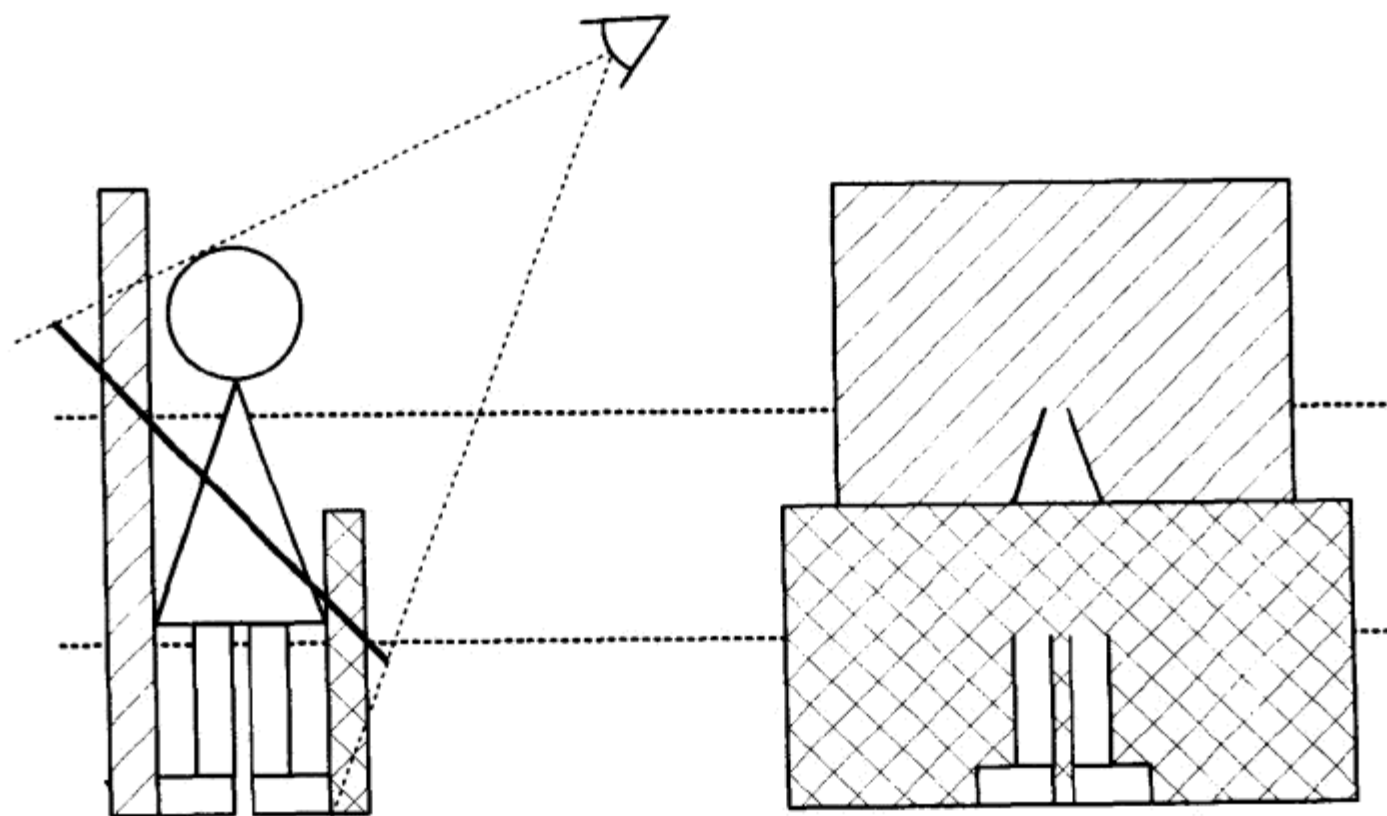


图 5.7.1 侧视以及渲染显示 Z 缓冲问题的广告牌替换体的视图

## 3. 矩形体 ( Cuboid )

下一个近似是包围盒。它绘制物体的物体空间的包围盒, 而不是绘制四边形, 让替换纹理投影到它上面。因为包围盒是真实的 3D 物体, 所以它的 Z 缓冲属性可以比那些与屏幕对齐的四边形得到更好的控制。特别是, 它的 Z 缓冲属性独立于摄像机位置, 而只依赖于物体位置和方向性, 应用时需要考虑的事情较少。

包围盒实际上工作相当良好。日常生活中的大部分物体的边界体非常分明, 在包围盒间相交的物体很少。注意, 仅绘出包围盒的前面和后面。可采用任何一种方式, 不过, 我推荐绘制前端, 因为许多物体充满它们的包围盒, 因此使用前端得到一种视差 (parallax error) 转

换，类似于摄像机移动时的真实物体。

虽然包围盒比四边形更复杂，但也不是太复杂。因为一般的场景有成百上千的三角形，12 个 tri 不会突破三角形预算。与单个四边形相比，Z 缓冲稳定性的增加值得我们付出努力。

#### 4. 包围物体

存在大量包围盒不能很好近似、而且使用它们可以导致不必要 Z 冲突的物体。

选择替换体形状另一个要考虑的因素是视差 (parallax error)。大部分 3D 场景拥有大量静态元素，但摄像机仍然要移动通过该场景。当摄像机开始移动时，一个上面绘制了一些景物图片的包围盒，不可能长时间保持看上去像这些景物。虽然增加替换体更新的速率有所帮助，但它只是加快像素和三角形的绘制速率，好处并不大。眼睛在移动时，可以很快注意到视差。

使用接近于真实物体形状的替换物体，虽然仍然通过非常低的三角形计数，但能得到很好地改进包围盒。形状上的主要限制是：它需要完全装入这个物体；否则，替换体的图像可能在屏幕上比被使用的替换体大，而且图像的边缘不会画出来。另一个限制是，物体必须总是凸的，避免出现自我排序问题，因为替换体纹理通过 alpha 混合来绘制。

#### 5. 图像变形

替换物体的一个问题是，它既必须是凸的，而且要比它所表示的东西大。前一个要求用来避免像素被渲染两次，后一个要求避免某些像素根本不会被渲染。不过，这种必要条件意味着，对象旋转时（或者摄像机移动时），视差不会正确，因为替换物体限制到比它表示的图像大。对于完全凸的源物体，使用更复杂的三角形替换物体，有助于避免这种现象的发生。不过，对于非凸的源物体，不会得到明显的改进，替换物体必须保持凸的，而且不允许有三角形与凹下物体匹配。

处理以上问题的另一种方式是在每一帧每一个顶点处移动纹理坐标。涉及的三角形数量相当少，因此每个顶点处稍多一些的工作不会对性能带来太多的影响。原理相当简单，即算出从某个角度观察时，每个替换体的顶点位于真实物体（以及替换纹理图像）的哪个地方。随着视角的改变，顶点所在的贴图像素也将改变。因此，对于新的视角，跟踪从观察者通过替换体顶点到原始物体的直线。然后，算出物体的这一部分初始绘制到的贴图像素，并设置这个顶点的相应 UV 坐标。

#### 6. 进一步的改进

以上方案开销大，且需要一定技巧来实现。它还有许多问题，如顶点光线跟踪在物体上留下视觉上不太明显的“裂缝”，而且因为替换物体中低密度的顶点在图像上会产生大的变形。另一个问题是，如何处理替换物体边缘处，根本没有贴图到真实物体任何部分的顶点，这些顶点将发生什么情况呢？从物体的可见边到这些顶点可能需要某种类型的插值。应用程序不希望在运行时完成这类工作，尽管涉及的顶点不多。

我们在实际中发现的简单得多、而且也能很好发挥的方法是给每个顶点赋予一个“视差因子”，即观察者每移动一度的图像质素数量。这个因子通常由人工来调整，而且它在运行时容易决定顶点的质素坐标。对于一个替换物体顶点，只进行一次因子设置，人工设置每个物

体的8~10个顶点不会花很长时间。

也可以对大部分真实的物体自动生成这些视差值,找出离每个边界物体顶点最近的真实物体顶点。这个距离与所要求的视差因子成比例。实际值取决于因子应用到UV值,又取决于纹理如何贴图到替换物体。对于更复杂的物体,如高度不规则、长而尖的、有空洞的或多层(multilevel)物体,最好办法仍然是人工设置这些因子。这个方法查找最近顶点,但对于这些物体,可能使用某类平均值更合适。

即使是这个简单的方法,也可能不能实现,因为物体是动态生成的(例如,通过动画)。在这种情况下,使用包围立方体,并在内部采用椭圆物体可以很好地工作。它比根本没有贴图坐标改变的情况工作得要好。

### 5.7.3 更新试探法

在每一帧上,对于每个替换体,需要决定是否更新它。作出此决定需要考虑大量因素。对每个因子,计算某种屏幕空间误差估计值,再加合所有因子。如果此和超过全局因子(可能在一定程度上为静态、特定的,或动态维持一定帧频的),则更新替换体。

#### 1. 动画

动画改变物体的外观,有些时候,误差会越变越大。计算渲染替换体所用的动画帧与未替换物体时使用的动画帧之间的差可以衡量这种误差。进行这种衡量需要物体在每一帧产生动画,即使不需要更新时也是如此。渲染物体的开销可能相当大,所以避免完全动画制作是一种明智选择。我发现的最简单途径是,对每个动画进行预处理,查找动画期间顶点移动的最大距离,除以动画的长度(用时间来衡量),于是得到衡量该动画的“每秒最大误差”。按蛮力方式很容易完成以上工作,而且因为是预处理,这是完全合理的。

注意,人眼绝对擅长于从有几个像素的图像中抽取出具有人类特点的运动。替换处理这个运动,大大减少它的有效帧频,变化可能非常显著,即使变化级别不大的情况下也是如此。一个不错的办法是,对于这些动画施加一个额外的偏移,与简单数学屏幕空间误差相比,它可以突出这些动画。它有效地消除对远处物体轻微动画的替换处理。

#### 2. 光照

如果光照在物体上变化幅度大,则需要更新。因为光照系统在3D引擎间变化极大,所以需要引擎针对性很强的程序来决定等价屏幕空间误差是什么。使用6个卡笛尔法线光照物体中央的某个点,而且比较当前条件与创建替换体时的条件的RGB差异,可以很好地了解光照的变化。乘以物体大小,并除以离摄像机的距离,可得到近似屏幕空间误差。

#### 3. 视角

改变视角大概是决定替换体更新的最重要的因子了。注意,重要的是从摄像机到物体空间中物体的向量。它将改变物体何时旋转,也改变摄像机何时移动,这两个信息都很重要。摄像机的观察方向不重要(除非使用了无数视场),当摄像机旋转时,外观不会改变很多,只在移动时发生改变。



#### 4. 摄像机距离

与物体到摄像机的方向一样，物体到摄像机之间的距离也同样重要。虽然这不改变物体的实际外观，但当摄像机直接移向物体时，替换体纹理将逐渐放大。一段时间后，可以明显看到，这只是一幅通过双线性滤波放大的图像，而不是真实的多边形物体，因此需要更新。这个更新将渲染到较大的纹理，得到更多的细节。

#### 5. 特定游戏的试探法

许多游戏也具有特定的试探法，可用于调整这些更新率。对 FPS 常用的一种是，接近视图中心的物体通常为玩家正盯着的物体。通过降低可接受的屏幕误差，应当稍微提高它们的更新速率。

对于鼠标驱动的游戏，如上帝游戏（god game）和 RTS 游戏，基于相同的原因，可以在鼠标光标之下对物体作类似调整。

在“风景”与“重要”物体间也可以加以区别。那些纯粹存在于风景中，用以创造一种气氛，但在玩游戏一般不涉及的物体，可以分配一个相当大的屏幕误差。玩家不会非常严密地审视它们，玩家的注意力不在此。因此，在玩家注意到这些误差之前，这些物体可以偏离得更远一些。

### 5.7.4 效率

---

在大部分卡上对效率的主要影响是渲染目标的改变。这会导致多个内部缓存和状态刷新，而且在某些卡上，将导致整个渲染管线的刷新。主要效率目标是最小化这些改变。完成此任务的最佳途径是等到当前场景的结束，批处理所需要替换纹理进行的更新，而不是在需要时马上进行。使用大型渲染目标，在场景的结尾，挑选出最空的渲染目标，并渲染多替换图像到该纹理的子区域。

### 5.7.5 预测

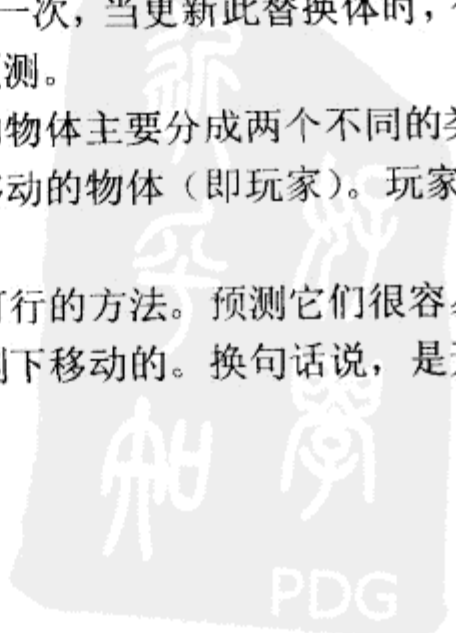
---

更新替换体时，物体的当前状态被渲染到纹理。然后在几个帧中淡入，接着，在被新的渲染替换前，在接下来的几个帧中保留一段时间。这意味着，在屏幕上所看到的总是过时的。

可以通过对替换体的状态进行预测来进行改进。其思路是，预测替换体在一半的生存期间将是什么样。如果某个物体目前是每 6 帧更新一次，当更新此替换体时，替换体的状态（方向、位置、光照、动画等）应当提前 3 帧进行预测。

对于第一人称视角射击游戏，游戏世界中的物体主要分成两个不同的类别，几乎从不移动的物体（墙、家具、其他点缀物）和无规律移动的物体（即玩家）。玩家的移动很难预测，而且试图替换玩家很可能是浪费时间。

另一方面，替换场景和家具尚不是一种更可行的方法。预测它们很容易，因为它们几乎从不移动。当它们移动时，通常是在玩家的控制下移动的。换句话说，是无规律的。最容易的事就是在移动过程中禁止替换。





对于上帝游戏 (god game) 和实时策略 (RTS) 游戏, 问题相似, 但摄像机的移动大不一样。它通常是一种鸟瞰视角, 而且大部分时间要么是静态 (对单元发出命令时), 要么以恒定速度在地图上移动, 到达不同的区域。小型、无规律移动很少, 而这些正是难以预测的, 所以对我们有利。在这些游戏中, 预测物体的移动可能非常有用, 因为它们大部分是由 AI 控制的。许多时间, 物体或为静态或以直线行进到目的地, 这两种情况都容易预测。不过, 摄像机移动和物体移动都可能突然改变, 而且当此情况发生时, 最好是立即标记该替换体进行更新, 甚至可以同时临时禁止替换处理。

### 5.7.6 小结

当绘制包含大量静态物体的场景时, 替换处理可以发挥作用。原始三角形计数会压垮任何试图以顶级细节来渲染它们的总线和图形设备。渐进网格方法只能在一定程度上减轻工作负荷。按这种方式, 极难减少纹理变化和动画。

对于离摄像机一定距离的静态物体, 替换体最有效。在游戏中引入这类点缀体大大提高了视觉质量。尤其是, 因为每个物体仍然是一个真实的独立 3D 物体, 在玩家需要的时候, 可以与它们进行交互。在大量物体中, 也可以让关键物体“消隐在清晰视线中”。通过现有技术以及限制场景中可见的物体数量, 这是很难实现的。

甚至使用包围盒和某些简单数学运算的实现, 对于游戏中的偶然物体, 也会产生好效果。这些偶然在游戏当前状态下不可见, 但可产生更加真实的场景。



## 5.8 硬件加速过程纹理动画中的运算

Greg James, Nvidia Corporation  
gjames@nvidia.com

**消**费级别的 3D 加速卡近年来取得了很大改进。当今最先进的芯片每一遍采样 4 个纹理，并提供强大的纹理寻址和像素处理运算。其中有依赖性的纹理读取，即从一个纹理采样的颜色值将用于影响后续纹理读取的坐标。依赖性纹理读取结合该功能渲染到颜色缓冲，并使用该缓冲作为后面渲染中的源纹理，这样就可以完全在图形处理器上生成感兴趣的纹理和纹理动画。采用接近于每秒 10 亿的纹理（贴图）像素（texel）采样率，这些过程纹理效果处理运行得非常快，而且对于实时 3D 加速场景很实用。

在计算机图形学的早期，我们就开始使用程序纹理生成的技术。模拟自然现象和生成复杂自然出现模式有多种不同的算法[Ebert98]。过程纹理动画擅长获得各种不同效果；而使用预渲染或“固定”动画只需内存和存储空间的少部分。容纳 2 个或 3 个全帧的内存通常就是生成无穷非重复动画所需要的全部内存。用户输入可以自由应用于这些动画上，而且这种交互性可以生成更丰富、更动人的虚拟场景。

本篇技术文章讲解了生成过程动画的少数基本运算，并应用于具体的实例中来模拟火焰、烟雾、水，或者执行图像处理。通过当今的消费级硬件，甚至可以完全在 3D 加速卡之内运行复杂的单元式自动程序，并将生成的动画应用于不同的效果中。

### 5.8.1 硬件运算

我们将着重讨论两种渲染运算，它们在生成过程纹理效果时用到。第一个操作是用相邻的贴图像素进行 4 重多重贴图采样，第二个运算是 2D 依赖性绿-蓝纹理寻址。两个最常见的 API（OpenGL 和微软的 DirectX8）支持这些运算。对于 4 路多重采样（four-way multisampling），我们将使用一个顶点程序，当引用顶点位置和纹理坐标时，将这个程序载入图形处理器中。这些顶点程序为 DirectX 8 的“顶点着色器”[DX8'00]，现在是 OpenGL 的 NVIDIA NV\_vertex\_program 扩展[NVExt2001]。对于我们所举例子，在为像素引擎提供输入的 4 个纹理阶段的每个阶段，顶点程序为邻近纹理像素建立采样建立适当坐标。像素引擎将针对不同效果综合这些样品，也可通过依赖性纹理寻址运算进一步操纵这些样品。像素引擎也是可编程的，它通过 DirectX 8 的“像素着色器”（Pixel Shader）和 NVIDIA 的

NV\_texture\_shader 扩展提供接口。首先，我们先查看顶点处理操作的情况。

针对模糊、卷绕以及物理过程的相邻采样

许多过程纹理的算法依赖于采样一个纹理像素的相邻像素，并过滤或模糊这些样本来创建一个新的颜色值。我们可以使用 4 重的多重贴图纹理采样把一张源贴图渲染到一个相同分辨率的颜色缓冲中，用来完成此源贴图的相邻采样。这张源贴图被选入到所有贴图单元中<sup>1</sup>，并用 Vertex Shader 程序生成 4 组独立的贴图坐标<sup>2</sup>，每组坐标偏移到一个相邻纹理像素的位置。使用从 0.0 到 1.0 的贴图坐标（可以精确覆盖整个渲染目标）来渲染一个四边形，并设置每组贴图坐标偏移量为 0，每个目标纹理像素将从源贴图采样同一个像素四遍。这样我们可以得到此源贴图的精确拷贝。通过矢量偏移每个贴图到相邻纹理像素，每个目的像素就从源纹理中相应像素相邻的 4 个纹理像素采样。我们就可以在像素引擎中卷绕这些采样像素（结合各种比例因子）。

如果我们使用点采样（point sample），4 重多重贴图（multitexture）的硬件最多给我们每渲染遍次（per pass）4 个相邻点采样。如果我们打开双线性过滤贴图模式，每个偏移采样会利用 4 个纹理像素，这样我们就可以在每渲染遍次采样 16 个相邻纹理模式<sup>3</sup>。在 2×2 双线性过滤中的纹理像素权重取决于精确贴图坐标的放置。例如，我们可以把 4 个双线性采样点正确地放置于两两相邻纹理像素之间，来获取一个纹理像素周围所有相邻点（8 个）的平均值。图 5.8.1b 显示了此操作<sup>4</sup>。程序清单 5.8.1 是在 4 重多重贴图硬件上进行相邻像素采样的代码。这些采点从源贴图获得，用于加和平均得到 X 点的像素值。我们可以使用 Vertex Shader 中的地址寄存器 a0.x，选择不同的采样位置（在这里指 A 或 B）。采样方案的索引值在渲染之前被调入到 Vertex Shader 的常量寄存器内。

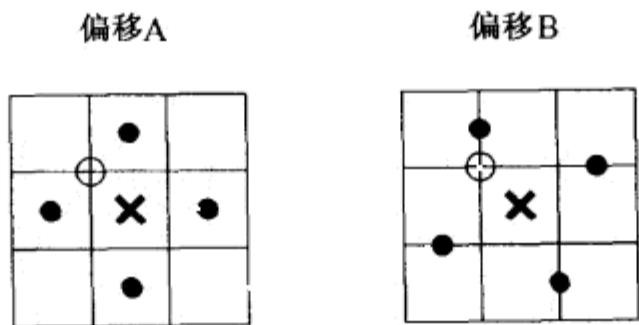


图 5.8.1 与程序清单 5.8.1 对应的贴图纹理像素采样。显示了标记为×的纹理像素被渲染时的采样模式。在每个贴图单元中（T0-3）求得的 Vertex Shader 程序的贴图坐标用点标记。圆心标记出当纹理坐标偏移为(0,0) 时采样点落入的地方，说明了 s\_off 和 t\_off 为半个像素偏移量。根据程序 5.8.1，每个纹理像素尺寸为(s1, t1)

1. 注：multitexture 的所有贴图 unit。  
2. 注：这里缺省认为 3D 显卡的多重贴图单元数为 4。  
3. 注：双线性过滤贴图是使用一个纹理像素周围的 4 个最近像素颜色，与此像素值进行带权值的加和后平均，得到最终在屏幕上  
的贴图纹理颜色。使用 4 重多重贴图的硬件（如 NV GEFORCE3），可以对一个纹理像素进行 4 次这样的操作，即每重贴图采样 4  
个相邻点，4 重多重贴图就采样了 16 个点。  
4. 图中每个黑点，代表多重贴图的一重。一共 4 重多重贴图（multitexture），所以有 4 个黑点。每个黑点在双线性（Bilinear）过  
滤贴图时，一个黑点的纹理采样要使用这个黑点周围最近 4 个相邻像素进行加和后平均。按图中所示，每个黑点精确地落在两个  
相邻像素中间，每个黑点周围最近的 4 个相邻像素恰恰两两重合，所以 4 个黑点的双线性过滤采样正好把中间像素周围的 8 个  
像素都采样到了。

### 程序清单 5.8.1 用于采样每个纹理像素邻居的代码与顶点程序

`RenderFullCoverageQuad()`通过每根轴上 0.0~1.0 的输入纹理坐标（准确覆盖渲染目标）渲染单个四边形。这些坐标按 4 种方式偏移得到 4 个输出 `oT[0-3]` 纹理坐标，因此，渲染每个像素时，它利用这些邻居纹理像素生成结果。所有大写变量通过 `#defines` 指令定义为顶点着色器常量内存的索引。

```
float s1 = 1.0f / texture_resolution_x; // 一个纹理像素的宽
float t1 = 1.0f / texture_resolution_y; // 一个纹理像素的高
float s_off = s1 / 2.0f; // 采样纹理像素中央
float t_off = t1 / 2.0f;

//4 个最近相邻点的 s,t,r,q 偏移量（双线性或点采样）
float offset_a1[4] = { -s1 + s_off, 0.0f + t_off, 0.0f, 0.0f};
float offset_a2[4] = { s1 + s_off, 0.0f + t_off, 0.0f, 0.0f};
float offset_a3[4] = { 0.0f + s_off, t1 + t_off, 0.0f, 0.0f};
float offset_a4[4] = { 0.0f + s_off, -t1 + t_off, 0.0f, 0.0f};

//8 个相邻点的 s,t,r,q 偏移量（双线性采样）
float offset_b1[4] = { s1/2.0f + s_off, t1 + t_off, 0.0f, 0.0f};
float offset_b2[4] = { -s1 + s_off, t1/2.0f + t_off, 0.0f, 0.0f};
float offset_b3[4] = { -s1/2.0f + s_off, -t1 + t_off, 0.0f, 0.0f};
float offset_b4[4] = { s1 + s_off, -t1/2.0f + t_off, 0.0f, 0.0f};

SetVShaderConstants((T0_BASE .. T3_BASE) + SET_A, offset_a1 .. offset_a4);
SetVShaderConstants((T0_BASE .. T3_BASE) + SET_B, offset_b1 .. offset_b4);
SetVShaderConstants( OFFSET_TO_USE, use_a ? SET_A : SET_B );
RenderFullCoverageQuad();
```

#### 顶点程序

```
; v0 = vertex position
; v1 = vertex texture coordinate

; Transform vertex position to clip space. 4-vec * 4x4-matrix
dp4 oPos.x, v0, c[ WORLDVIEWPROJ_0 ]
dp4 oPos.y, v0, c[ WORLDVIEWPROJ_1 ]
dp4 oPos.z, v0, c[ WORLDVIEWPROJ_2 ]
dp4 oPos.w, v0, c[ WORLDVIEWPROJ_3 ]

; Read which set of offsets to use - set A or B
mov a0.x, c[OFFSET_TO_USE].x

; Write S,T,R,Q coordinates to all four texture stages, offsetting
; each by either offset_a(1-4) or offset_b(1-4)

add oT0, v1, c[ a0.x + T0_BASE ]
add oT1, v1, c[ a0.x + T1_BASE ]
add oT2, v1, c[ a0.x + T2_BASE ]
add oT3, v1, c[ a0.x + T3_BASE ]
```

请注意，把一个与源纹理相同分辨率的四边形渲染到缓冲时，一个偏移量为 (0,0) 的贴

图坐标将从一个纹理像素的左上采样（最低的坐标点）。要从纹理像素中央采样，必须在偏移上添加一半的纹理像素宽度和高度，或者在每根轴上移动四边形半个像素。程序清单 5.8.1 选择添加一个纹理像素宽度和高度。当使用双线性过滤时，有必要理解这些半纹理像素偏移。没有它，双线性样品将获取错误的相邻像素，并生成大不一样的结果。有必要测试并知道纹理样品的位置，让程序纹理算法按要求进行工作。本文后面描述的“Conway 的生命游戏”就是一个很好的测试实例。

得到的 4 个贴图采样结果<sup>1</sup>可以在可编程的像素引擎<sup>2</sup>中混合起来，得到各种效果。把结果图像当作新的源纹理，并再次应用邻居采样创建后续帧，可以得到各种不同的有趣纹理动画。如果平均 4 个样品，可得到源图像的模糊效果。对前述偏移引入一个附加的(s,t)滚动量，可以在连续多帧中模糊和滚动源图像。模糊和滚动时，可以抖动用于每一帧的滚动向量，并且在样品颜色上乘以不同的 RGBA 值，渐变或更改颜色。如果应用稳定的输入来模糊和上滚，首先渲染明亮源像素或每一帧底部的“灰烬”，则得到图 5.8.2 所示的火焰和烟雾效果。只使用两个 128×128 的 32 bit 纹理，在现代图形卡上，该效果以每秒 500 帧的速度不重复地延续。因为我们不能同时渲染到纹理并将它作为源，所以，必须使用两个纹理，而且在它们之间来回往复。一个作为前一帧源，另一个作为当前帧目标。

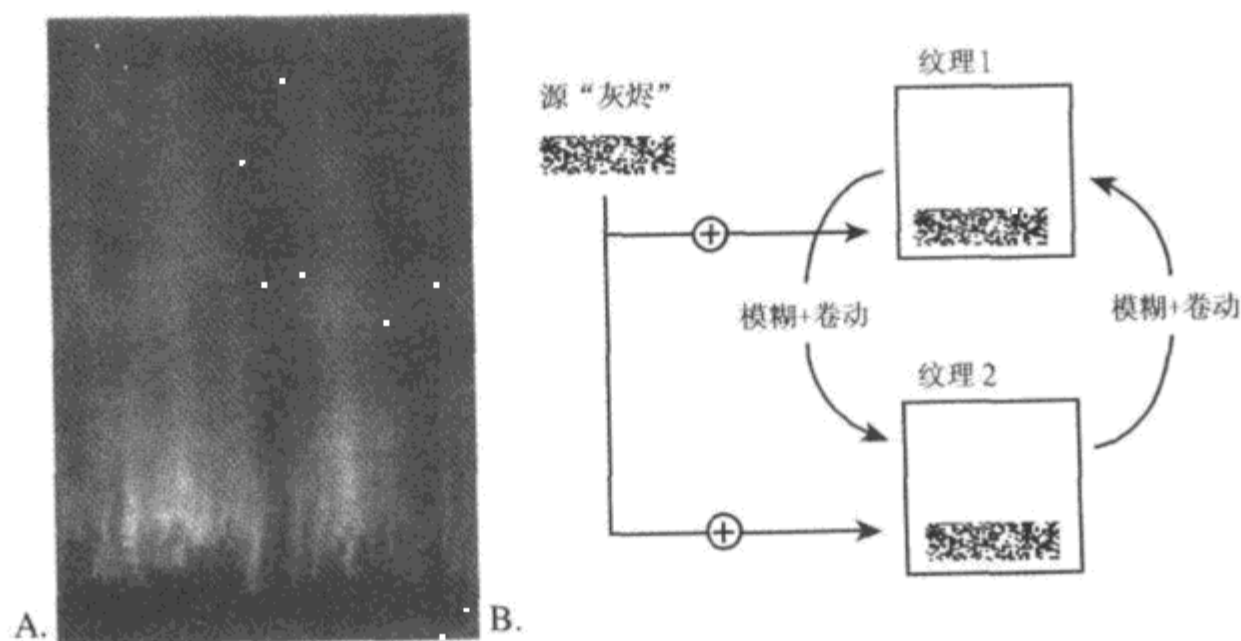


图 5.8.2 使用程序清单 5.8.1 的偏移 A 以及滚动偏移得到的火焰和烟雾动画。在模糊和上滚动之前，底部的明亮“灰烬”渲染到每一帧的纹理。图 B 显示了渲染运算的过程

不用平均邻居样品来模糊，我们可以计算样品之间的差异，并取幂，以提到高的对比度。使用程序清单 5.8.1 的偏移 A，并减少偏移的幅度，对中央质素进行部分采样，可以在硬件 Pixel Shader 程序中执行边缘检测，请参考 Matthias Wloka 撰写的论文[Wloka2000]（如图 5.8.3 所示）。另一种奇异的“毛玻璃”效果可通过模糊源并差分模糊与源图像得来。图 5.8.3d 显示了  $(src \text{ DIFF } (src \text{ DIFF } (BLUR(src))))$  的结果，其中 DIFF 为 RGB 颜色差的绝对值。

1. 注：多重贴图的每个贴图单元得到一个采样结果，因为这里只用 4 重多重贴图硬件为例，所以是得到 4 个采样结果。  
2. 注：这里指 Pixel Shader 或者 OpenGL 中的 ARB\_fragment\_shader，NV 扩展是 NV\_register\_combiners，ATI 扩展为 ATI\_fragment\_shader。



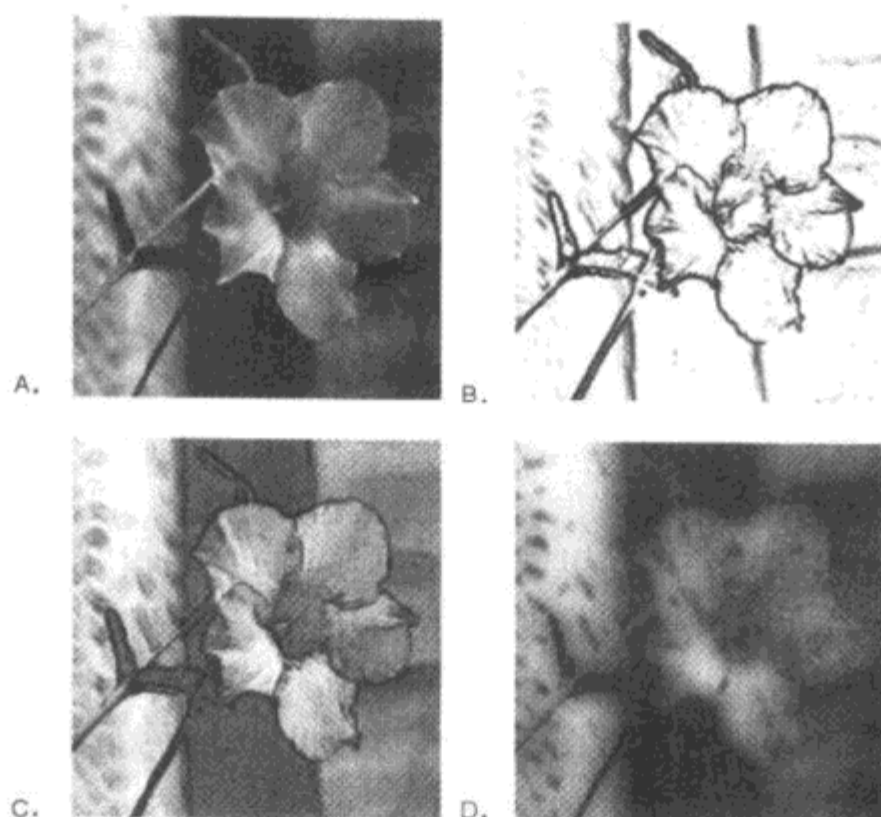


图 5.8.3 边缘检测与图像处理。A)原图；B)在可编程像素硬件中边缘检测的结果；C)A 和 B 50%的混合；D)原图减去原图与原图模糊后的差

邻居取样和差分也可以用于通过连续渲染到纹理运算实现基于高程场的水物理变化。《游戏编程精粹 1》中的文章[Gomez00]介绍了通过 CPU 上的顶点几何处理，用于高程场水变化的相关算法。我们可以将它扩展成一系列的渲染运算，使用纹理来表示水高度、速度以及强度。每个纹理像素取代 Gomez 实现程序中的一个顶点。在这里，我们不是采样 CPU 上的邻居顶点，而是采样图形处理器上的相邻纹理像素。我们的实现用到 6 个纹理。两个纹理表示水高度，作为灰度颜色值（高度地图），一个纹理表示当前时长的高度，一个纹理表示前一个时长的高度。与之相似，用两个纹理表示水速，两个纹理用于累加作用于每个纹理像素上的最近邻居的强度。图 5.8.4 演示了使用这一技术的动画中的 4 个帧。尽管使用了 6 个  $256 \times 256$  的 32 bit 纹理，而且每一时长有 4 渲染遍次，但动画仍然保持超过每秒 450 个时长的速率。图 5.8.5 显示了用于生成一个时长中的纹理的过程。

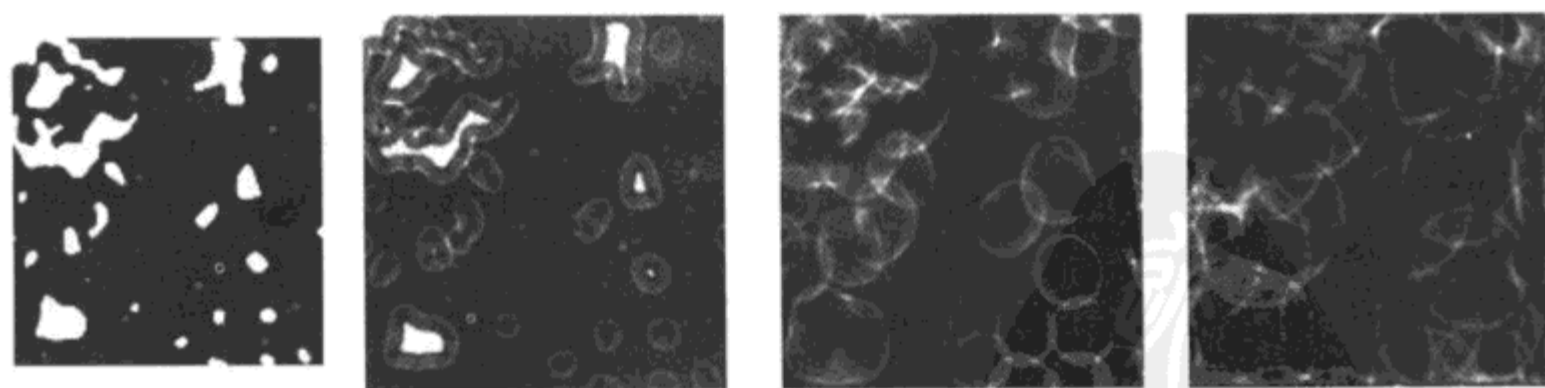


图 5.8.4 3D 加速卡处理像素的过程中，初始动画以及基于高度的水动画的 3 个帧。在生成后续时长时使用了 6 个纹理，不过，这里只显示了输出高度纹理



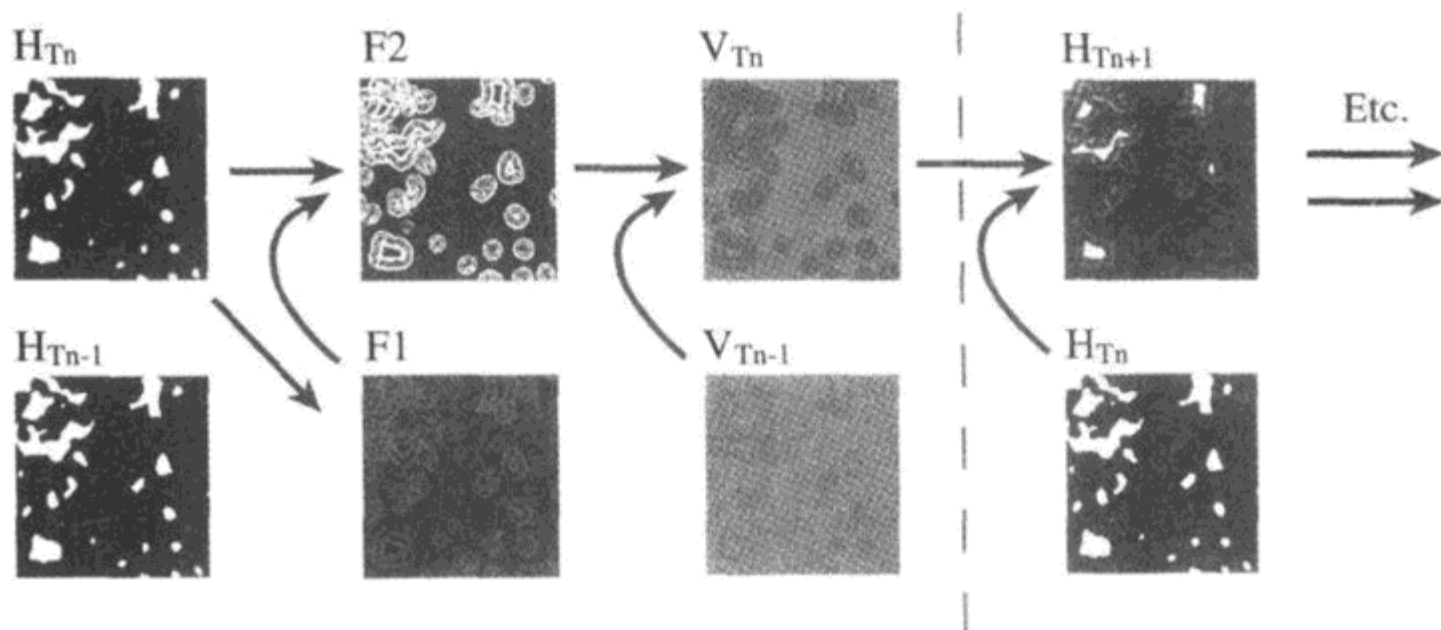


图 5.8.5 基于高度的水动画使用的 6 个状态纹理。 $H_{Tn}$  是最近纹理时长的高度。 $F1$  和  $F2$  用于累加作用于每个纹理像素高度的强度。得到的  $F2$  应用到前一个时长的速度  $V_{Tn-1}$ ，得到的速度  $V_{Tn}$  应用于高度来创建下一时长  $H_{Tn+1}$  的高程场。图中没有显示乘到纹理像素值上的“质量”和“时间”的比例因子

此例中使用的硬件上，每个分量有符号颜色值内部运算于 9 bit 精度，这将导致不准确结果。低位容易丢失，在物理过程中产生舍入误差，导致系统衰减到 0，或者增长到饱和。最近邻居采样也可能在以噪音（noise）形式出现的纹理像素值之间产生高频振动。通过模糊和添加恢复强度，将所有的高度平缓降到中间范围值，衰减系统后可以减轻这些问题的严重性。模糊是使用双线性过滤采样偏移大于一个纹理像素宽或高的相邻点来完成的。

动画的灰度高度图是很有用的，但我们可以进一步使用最近相邻点采样，利用一张灰度图来创建一张 RGB 法线图。这种法线图可以用于逼真的点积凹凸贴图[Moller99]。法线图的创建可以在图形处理器上，使用“渲染到纹理”技术一次渲染完成，这样做避免了 CPU 开销、贴图内存拷贝和图形管道的延迟。用这种方法在硬件上创建法线图是一种优雅的手段，用来更新和动画表面的细节。灰度高度特征如子弹孔，裂缝等可以在游戏运行时渲染到高度图中，再由灰度高度图转换为法线图。转换的速度与前面讨论的火焰和烟雾效果的效率相当。程序 5.8.2 显示了一个 PixelShader 程序，在 4 重多重贴图的硬件上可以用一遍渲染把一张灰度高度图转换为 RGB 法线图。它使用近似计算方法把 RGB 结果值归一化。其结果对于光照和反射计算是非常理想的。使用两遍渲染与依赖型纹理读取操作可以用来产生精确化的归一化 RGB 向量。

使用程序清单 5.8.1 中的纹理像素偏移 A。

程序清单 5.8.2 根据灰度高度图，在硬件中创建 RGB 普通图的代码。在所有 4 个纹理阶段选用此灰度图像，同时使用程序清单 5.8.1 中的偏移 A

### (1) 像素程序

```
// 从一张输入的灰度高度图创建一张 RGB 法线图
// 没有使用配对的 RGB 和 Alpha 指令
// 法线图参数化范围 (0,1)，所以 0.5=轴的零分量
// 红色=正 s 轴方向
```

```

// 绿色=正 t 轴方向
// 蓝色=正 r 轴方向 (垂直纸面向外)

// Declare pixel shader version
ps.1.1

def c5, 1.0, 0.0, 0.0, 1.0 // s 轴的红色掩码
def c6, 0.0, 1.0, 0.0, 1.0 // t 轴的绿色掩码
def c4, 0.0, 0.0, 1.0, 1.0 // r 轴的蓝色掩码(blue = up out of texture)
def c2, 0.5, 0.5, 0.0, 0.0 // 红与绿的 0.5 偏移
def c1, 1.0, 1.0, 0.0, 0.0 // 红与绿的颜色掩码

; get colors from all 4 texture stages
; t0 = -s, 0
; t1 = +s, 0
; t2 = 0, +t
; t3 = 0, -t

// 4 个贴图单元全部选择源灰度高度图
// 4 个贴图单元都进行采样
tex t0 // t0 = 位于坐标偏移 (-s, 0) 的 RGBA 像素纹理
tex t1 // t1 = (+s, 0)
tex t2 // t2 = ( 0, +t)
tex t3 // t3 = ( 0, -t)

sub_x4 r0, t0, t1 // r0 = (t0-t1)*4 = s 轴的高度斜率
// 使用_x4 来增加输入灰度图的对比度
mul t0, r0, c5 // t0 = r0 * 红色掩码 = 红色分量
// 使用 t0 作为临时变量

sub_x4 r1, t3, t2 // r1 = (t3-t2)*4 = t 轴的高度斜率
mad r0, r1, c6, t0 // r0 = r1.green + t0 = 在 red 和 green 中 s 和 t 的结果
mul t1, r0, r0 // t1 = s 和 t 分量的平方
// 使用 t1 作为临时变量

dp3_d2 r1, 1-t1, c1 // r1.rgb = (1 - s^2 + 1 - t^2 )/2
// 当 s 很小, (1-s^2) 为 sqrt(1-s^2) 的近似
add r0, r0, c2 // r0 = r0 + 0.5 red + 0.5 green
// 把符号值移位到 (0,1)

mad r0, r1, c4, r0 // RGB = (r+0, g+0, 0+blue )
// output = r0

```

## (2) 依赖性纹理寻址

我们已知道, 简单多采样运算可以成功生成不同的效果。还可以使用另一个更复杂, 而且更强大的纹理运算。通过依赖性纹理寻址, 我们可以在另一个纹理颜色的基础上, 从一个纹理取出纹理样品, 或者在重复纹理坐标和纹理颜色的基础上获取。本篇文章只介绍更简单依赖性纹理运算中的一种: 依赖性绿-蓝纹理寻址运算, 用 DX8 的 `texreg2gb` 指令表达, 或者用 OpenGL 的 `GL_DEPENDENT_GB_TEXTURE_2D_NV` 扩展来表达。

依赖性绿-蓝纹理寻址是一个简单的过程。对于被渲染的像素, 硬件获取源纹理在该像

素的绿色分量用作读取另一个纹理的 S 坐标（水平坐标），蓝色分量用作 T 坐标（垂直坐标）。如图 5.8.6 所示，它具有一个 3×3 的源纹理。

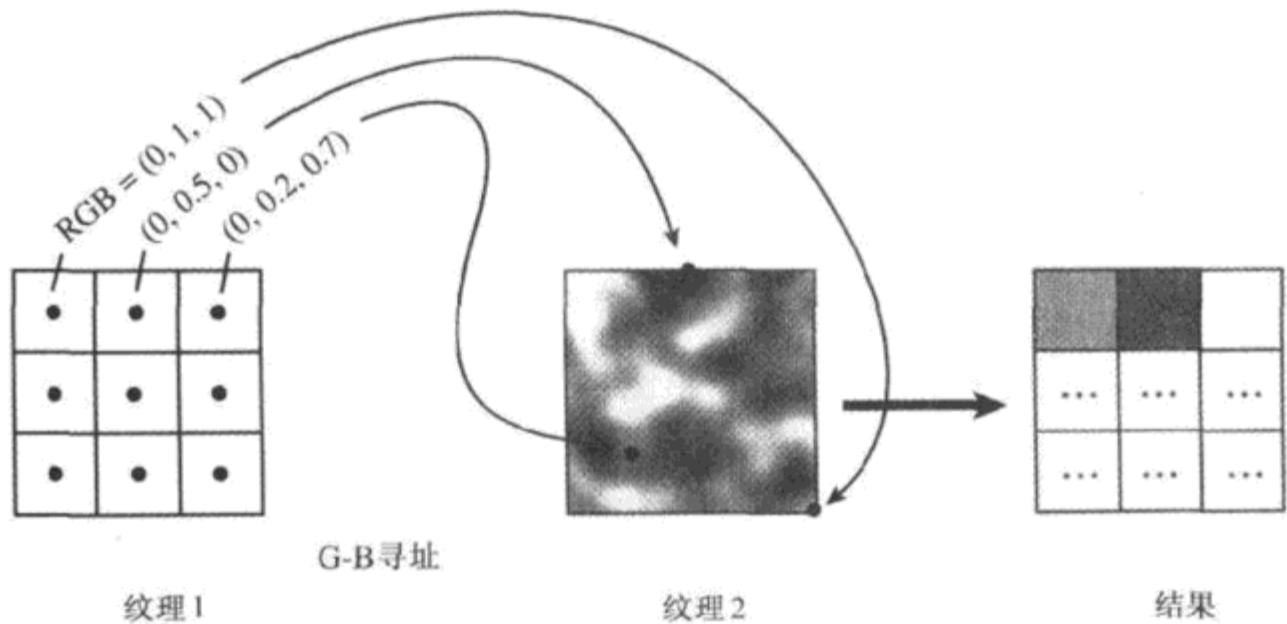


图 5.8.6 依赖性绿—蓝纹理寻址。纹理 1 源在指示点采样，纹理 1 的每个纹理像素决定从纹理 2（即“规则”纹理）采样的坐标。从纹理 2 得到的颜色作为结果的输出

纹理 2 相对于纹理 1 可以为任意大小，不过，如果纹理 1 包含 8 bit 的绿色和蓝色值，则对于纹理 2，任何分辨率大于 256 的纹理都无意义，因为通过粗糙的 8 bit 值不能访问更精细的分辨率。纹理 2 提供了纹理 1 输入的任意查找表，而且可以与渲染到纹理运算一起使用，完全在像素硬件内运行一些非常复杂的程序。纹理 1 称做“输入”纹理，纹理 2 称做“规则”纹理，因为纹理 2 决定输入色如何贴图到结果。

通过渲染到纹理的运算，可以使用一个或多个纹理贴图来保存中间逻辑结果。也可以在运行时基于前面结果或用户的输入生成“输入”和“规则”贴图。正如你所想，这里存在极大的灵活性，而我们甚至还没有考虑依赖性 alpha-red 查找或点积寻址运算操作。因为绿—蓝寻址操作取决于特殊的贴图颜色，这些颜色对于场景显示可能并不是最佳的。所以可以在颜色表中加入第二个依赖查找，这个方法很好而且容易实现。这第二个操作把过程纹理的结果映射到任意 RGBA 值，以用来在屏幕上显示。它创建了一个最终输出的纹理，而用户则从来看不到驱动计算的绿—蓝纹理。最后一个例子综合相邻采样与依赖性查找在图形处理器上运行“Conway 的生命游戏”。

(3) 硬件中“Conway 的生命游戏”

John Conway 的“生命游戏”[Gardner70]是一个流行的单元自动机程序（cellular automata program）。虽然初看之下，它似乎与一般的计算游戏关系不大，或者无多大用途，但这是工作所熟悉的领域，而且在图形芯片上运行游戏的基本运算适用于其他程序技术。游戏可以生成相关噪音的有趣模式，它们在驱动其他效果中 useful。例如，前述的火焰和烟雾效果中的灰烬。在您的平台上实现这个游戏是验证硬件正完成所期望功能的良好途径。游戏具有可识别循环模式，而且敏感地依赖于正确的样品位置，这使得它很容易判断何时采样正确。

在“Conway 的生命游戏”中，每个单元或以“on”开始，或以“off”开始，由白色或黑色质素的纹理贴图来分别表示它们。对于贴图上的每个单元，创建下一次生成的规则是：

- 如果单元为“on”，而且有 2 个或 3 个邻居为“on”，在下代中单元保持为“on”。
- 如果单元为“on”，而且少于 2 个，或大于 3 个邻居为“on”，则该单元在下代中为“off”。
- 如果单元为“off”，而且有 3 个邻居为“on”，则该单元在下代中为“on”。

游戏要求采样 8 个邻居和中心单元，并应用逻辑到采样结果上。OpenGL 的“红宝书”[Neider93]讲解了使用硬件模版缓存运行游戏的技术，但通过硬件依赖性纹理寻址功能，可以利用较少的渲染遍数，用更灵活的方式来实现。

我们的方式是创建渲染目标颜色缓冲，其分辨率与单元域相同。缓冲中的每个纹理像素以黑色开始，对于 8 个邻居纹理像素和中心纹理像素，我们累积性地从源单元域渲染颜色。源的每个邻居纹理像素乘以  $1/8$  绿<sup>1</sup>，并加合到目标上，中心纹理像素乘以全蓝<sup>2</sup>并加合。实际上，我们使用双线性采样和程序清单 5.8.1 中的偏移来一次采样两个相邻像素，并将样品乘以  $1/4$  绿。第二遍采样乘以蓝色的中心纹理像素。其结果是一个编码了每个源纹理像素相关游戏规则的（生存）条件，这些纹理像素与游戏的规则有关。每个单元的邻居的计数范围从 0 到 1（绿色），单元的“on”或“off”状态为 0 或 1（蓝色）。于是每个绿—蓝纹理像素用作依赖性绿—蓝寻址运算的输入。接着从一张  $8 \times 2$  的“规则”纹理读取像素。这张规则纹理决定下一代纹理像素的颜色（黑或白）。这个规则纹理在所有纹理像素中为黑色（除(2,1)、(3,1)和(3,0)为白色外）。这些白色像素为那些按游戏规则决定在下代中为“on”的像素。图 5.8.7 演示了一个处于初始源纹理、根据它进行渲染的绿—蓝条件、规则纹理以及生成的后续代。

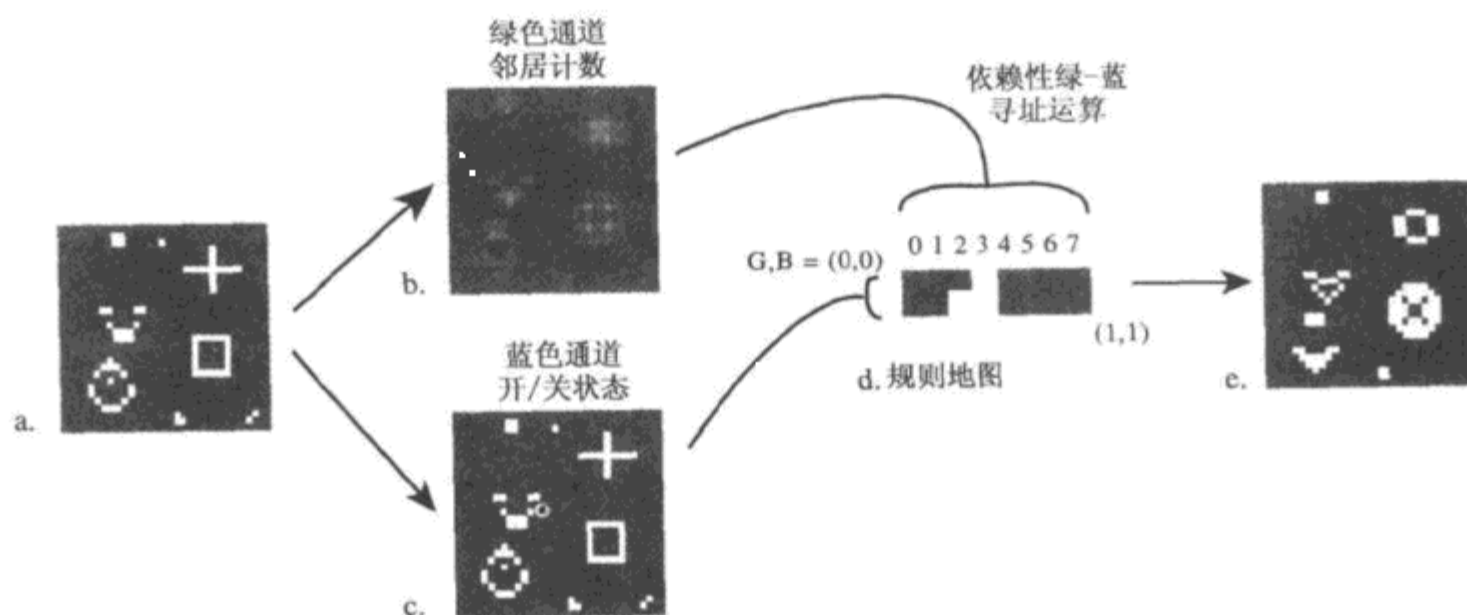


图 5.8.7 在“Conway 的生命游戏”中，生成下一代单元的步骤。a) 初始单元域。b) 为每个单元的激活邻居的条件纹理的绿色分量。c) 反映单元是否为“on”或“off”的条件蓝色分量。d) 编码游戏规则的  $8 \times 2$  像素纹理。e) 再次用于新输入 a 的单元的生成。

这个模拟取决于白色的源纹理像素以及绿—蓝的中间纹理像素。从源或中间纹理像素到

1. 注：绿色分量为 (0,1,0)。

2. 注：蓝色分量为 (0,0,1)。

这里不需要限制规则纹理为  $8 \times 2$  像素。可以装载任何大小的纹理,而且这个纹理可以编码非常复杂的规则,或任何 RGB 值的派生像素。请参考在线示例,以及下面所提的源码。

### 5.8.2 将来的工作

---

随着图形处理器每遍渲染绘制更多的纹理示例,当前硬件的 8 bit 颜色分量以及内部像素处理运算的精度将得到进一步提高。最终,清除颜色条带或其他人工渲染痕迹都变得很容易。下一代硬件必定需要更高的精度格式,每个颜色分量具有 16 bit 或 32 bit,非常适合于获得硬件加速的程序纹理效果。

在前面水模拟中提到的精度限制问题,采用 16 bit 浮点值几乎可以消除它们。纹理和像素计算中的高精度浮点值可以让液体流动模拟和基于网格的 Navier-Stokes 算法在图形硬件中运行[Witting99]。这些算法的 CPU 实现已经接近于特效和图形研究中的实现。渲染目标缓冲的精度增加也将允许进行涉及简单 Fresnel 积分、波传播以及衍射效果的硬件加速的高级光照计算[Feynman85][Stam99]。

单元自动机程序将导致具有不同程度相关性的大量有趣动画和模式产生。单元自动机程序具有难以控制引擎针对特定模式或效果的缺点,但在互联网上有大量现有解决程序以及方便浏览的程序。使用网格气体自动机(lattice gas automata)(一般是 2D 表面模拟)也可能提供有用的动画模式。

实时程序纹理和动画的将来令人鼓舞。可以利用大量技术和效果,而且它们必定具有给 3D 游戏带来增强真实性和多样性的功能。

### 5.8.3 致谢

---

特别感谢 Matthias Wloka, 他为多重采样以及硬件图像放置方面作了大量的努力。还要感谢 Nvidia 公司的同事,是他们提供了如此强大的图形硬件。

### 5.8.4 示例源码

---

针对前面例子的 DirectX 8 代码张贴在 Nvidia 的公共开发者网站上([www.nvidia.com/Developer/DX8](http://www.nvidia.com/Developer/DX8)),读者可以下载并进行试验。这些示例可以运行在 DX8 SDK 的参考光栅(reference rasterizer)中,或运行于那些支持 DX8 Vertex Shaders v1.1 和 Pixel Shaders v1.1 的硬件上。

### 5.8.5 参考文献

---

[Ebert98] Ebert, David S., et al, *Texturing and Modeling: A Procedural Approach*, Academic Press, 1998 [ISBN 0-12-228739-4].

[DX8'00] Microsoft Corporation, *DirectX8 SDK*, available online at <http://msdn.microsoft.com/directx/>, November, 2000.



[NVExt2001] Nvidia Corporation, "Nvidia OpenGL Extensions Specifications" (nvOpenGL Specs.pdf), available online at [www.nvidia.com/opengl/OpenGLSpecs](http://www.nvidia.com/opengl/OpenGLSpecs), March 2001.

[Wloka2000] Wloka, Matthias, "Filter Blitting," available online at [www.nvidia.com/Developer/DX8](http://www.nvidia.com/Developer/DX8), November 2000.

[Gomez2000] Gomez, Miguel, "Interactive Simulation of Water Surfaces," *Game Programming Gems*, Charles River Media Inc., 2000: pp. 187~194.

[Moller99] Moller, Tomas, and Haines, Eric, *Real-Time Rendering*, A K Peters, Ltd., 1999.

[Gardner70] Gardner, Martin, "Mathematical Games," *Scientific American*, vol. 223, no. 4, October 1970, pp. 120~123.

[Neider93] Neider, Jackie, et al, *OpenGL Programming Guide*, Addison-Wesley Publishing Co., 1993: pp. 407~409.

[Witting99] Witting, Patrick, "Computational Fluid Dynamics in a Traditional Animation Environment," *Computer Graphics Proceedings (SIGGRAPH 1999)* pp. 129~136.

[Feynman85] Feynman, Richard P., *QED: The Strange Theory of Light and Matter*, Princeton University Press, 1985, pp. 37~76.

Although he does not label it a "Fresnel integral," this is the name for the calculation he explains in Chapter 2. This powerful mathematics accounts for all phenomena in the propagation of light by a large sum of a few simple terms.

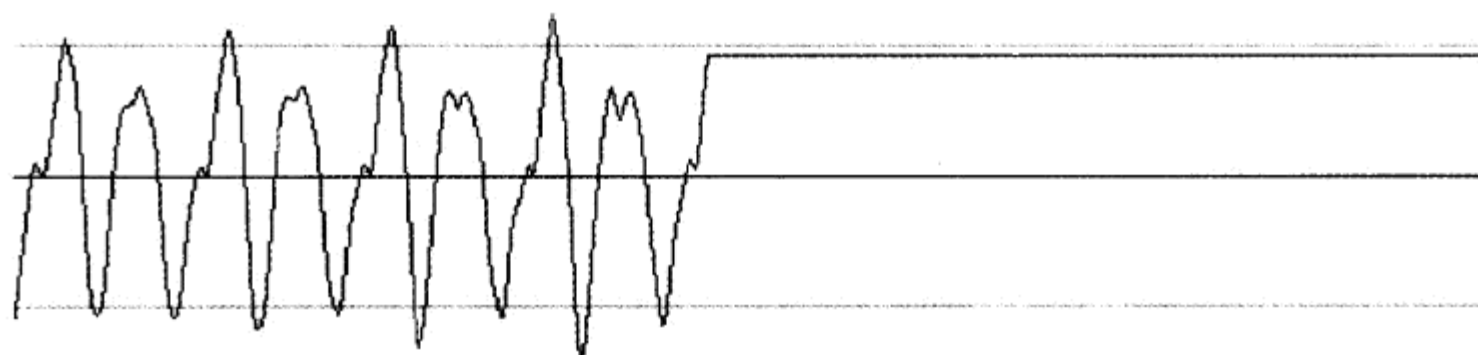
[Stam99] Stam, Joe, "Diffraction Shaders," *Computer Graphics Proceedings (SIGGRAPH 1999)* pp. 101~110.

Gomez, Miguel, "Implicit Euler Integration for Numerical Stability," *Game Programming Gems*, Charles River Media Inc., 2000: pp. 117~181.





## 音频编程





## 绪 论

James Boer, Lithtech Inc.

jimb@lithtech.com

欢迎来到游戏开发中数字音频的迷人世界！图形可以给生活带来交互的数字世界，但声音（voice）和音乐赋予游戏灵魂。没有音频的世界空虚乏味。毫不夸张地说，除了音频外，能如此强烈影响玩家情绪的方式很少。现在的游戏设计者刚开始发挥它的威力，恰似好莱坞在观众的心弦上演奏美妙的声音与动人的音乐。

在过去几年间，硬件与软件技术的发展给一度是游戏开发配角的音频处理带来了新的活力。游戏音频现在成为一个专门的研究领域，专门从事音频工作的程序员开始出现，这就像当前工业开始使用 AI、物理学、图形学专家一样。音频管线开发、数字过滤、3D 定位音频（3D positional audio）、闭塞和阻塞方法、API 设计以及交互音乐技术对于当今音频程序来说都是设计上的挑战。在本篇以及其他技术精粹文章中，我们将讨论以上主题以及其他相关内容。

在《游戏编程精粹 1》中，我们曾打算列入音频内容，但由于许多不同的原因，未如愿以偿。本书终于得以实现，但它仍然有一些难度。这也正反映了音频系统在游戏开发上为了作为独立领域而进行的奋斗。在过去，声音和音乐编程经常是在项目的最后几个月委托一些次要的程序员来完成的。现在，人们期望游戏音频像图形效果一样也具有震撼的效果和丰富的表现力。

但实现这些功能的实用解决方案仍然不太完善，一些资源和信息分散在互联网上，很少有专门的书籍来讲解游戏中的音频编程。现在，高级数字音频技术在游戏编程领域流行开来，本章技术精粹文章正是来源于一些著名的工业专业人员的总结。我希望这些文章对于解决读者的游戏音频编程难题大有帮助。



## 6.1 游戏音频设计模式

Scott Patterson

scottp@tonebyte.com

近几年来设计模式（design pattern）概念非常普及。实际上，设计模式已经被使用了很长一段时期。最近，我们对它们进行了更好地分类和识别。面向对象的编程语言通常和设计模式的实现联系在一起，但即使不是面向对象的编程语言也可以提供有效的对象模式实现。设计模式也在创建编码系统中起到指导作用。在本文中，我们也将把它们作为音频接口设计（audio interface design）的指导。这里假设我们是音频接口设计者，我们的客户是游戏程序员，我们要向客户提供一个方便、灵活和强大的音频接口。我们将简短总结一些设计模式以及如何将它们与音频接口设计结合在一起。

### 6.1.1 桥接（Bridge）

“将抽象部分与它的实现部分分离，使它们都可以独立地变化。”

声音标识符（Sound Identifiers）

一个将抽象部分与实现部分分离的有效途径是传递标识符，而不是传递类指针或引用等。标识符可以定义成数字或字符串。既然我们讨论的是音频，就要与每一个特定声音（voice）的标识符打交道。我们可以通过传递标识符来开始和结束一段声音。音频系统实际上如何开始和结束声音，与 API 客户端完全分离。

```
void StartSound( int nSoundId );  
void StopSound( int nSoundId );
```

同样，装载和访问声音的方式也被隐藏。我们可以定义一个类似的调用来加载和卸载特定的声音。

```
void LoadSound( int nSoundId );  
void UnloadSound( int nSoundId );
```

也许声音集合的另一套标识符系统会更有用。现在我们可以同时加载和卸载多个声音。

```
void LoadSoundCollection( int nSoundCollectionId );  
void UnloadSoundCollection( int nSoundCollectionId );
```

知道当前某段声音是否被装载可能有用。

```
bool IsSoundLoaded( int nSoundId );
```

如果尝试启动一段没有装载的声音，结果可能是无声或播放一段错误的声音。

### 6.1.2 外观 ( Façade )

---

“为子系统中的一组接口提供一个一致的接口。外观定义了一种更高层次的接口，使子系统更容易使用。”

#### 子系统控制 ( Subsystem Control )

当我们为游戏程序员编写音频 API 时，其目的是隐藏音频系统中所有繁琐的复杂性，但要满足游戏各方面的需要。带着这种目的编写 API，就像用外观设计模式设计类一样。音频系统的复杂性被游戏代码的复杂性掩盖了，而两个系统的连接就是 API。

有多种方法可以产生声音，但使用同样的接口来控制这些声音，使得音频子系统更加容易使用。设置和获取主音量的调用看起来十分简单，但实际上可能要更新内部的多个音频子系统。

```
float SetMasterVolume( float fVolume );  
float GetMasterVolume();
```

与软硬件合成、软件合成以及流代码有关的音频子系统可能被这样的调用全部更新，但这种复杂性被隐藏了。类似的复杂性在中止当前所有声音的调用中也可以被隐藏。

```
void StopAllSounds();
```

### 6.1.3 合成 ( Composite )

---

“将对象组合成树型结构以表示‘部分—整体’的层次关系。合成使客户用一致的方式对待单个和组合对象。”

#### 1. 引擎控制 ( Engine Control )

像汽车发动机一样的东西实际上可能由很多不同声音组成。这些声音可能由隆隆声、呜呜声、叮当声和噼啪声组成并同时响起。这些声音的音量、音调和其他参数可以映射成不同的游戏参数，如引擎类型、油门杆、功率大小、现行速度和许多其他参数。在这种情况下，我们希望针对所控制的引擎类型隐藏其复杂性并提供相应函数。

```
void StartEngine( CarInstance_t *pObject );  
void UpdateEngine( CarInstance_t *pObject );  
void StopEngine( CarInstance_t *pObject );
```

现在，内部音频代码可以解释已知的 CarInstance\_t 对象的状态和参数，并把它们转换成一段或多段声音的控制。无论是控制单个声音对象还是控制合成声音对象，客户可以用同样

的方式对待引擎控制。

## 2. 氛围控制 (Ambience Control)

对周围声音的控制也可以表现合成行为。如果我们要模拟一种环境，如丛林，可以随意地播放一些动物的声音。在这种情况下，在游戏中，我们可能没有一个叫 `Jungle_t` 的结构，但我们可能知道离丛林地区的距离以及动物的兴奋程度。

```
void StartJungle( float fDistance, float fActivity );  
void UpdateJungle( float fDistance, float fActivity );  
void StopJungle();
```

### 6.1.4 代理 (Proxy)

---

“为另一个对象提供一种代理或占位符，以控制对它的访问。”

#### 句柄 (Handle)

句柄为另一个对象提供一种占位符，以控制对它的访问。当我们启动一个特定的声音实例时，会希望持续地控制该实例。控制参数可以是任意从 3D 定位信息到声音的直接控制，如音量、音调、立体声左右均衡 (pan) 和效果。启动函数返回一个句柄，更新和中止函数使用该句柄来访问这个声音实例。

```
Handle_t StartHandledSound( int nSoundId , const ControlParams_t &cp );  
void UpdateHandledSound( Handle_t hSound, const ControlParams_t &cp );  
void StopHandledSound( Handle_t hSound );
```

关于句柄的讨论请参见[Bilas00]。

### 6.1.5 修饰器 (Decorator)

---

“动态地给对象添加额外的功能。修饰器为子类的扩展功能提供灵活的替代方法。”

#### 1. 用户数据 (User Data)

一种允许赋予对象动态功能和关联的方法是提供用户数据访问。如果我们提供一个用户数据域，让用户可以为一种声音的每个实例设置这个数据域，则可以帮助用户将其他功能和声音实例链接起来。我们可以将用户数据访问函数加到句柄化的声音界面上。

```
void SetHandledSoundUserData( Handle_t hSound, UserData_t UserData );  
UserData_t GetHandledSoundUserData( Handle_t hSound );
```

#### 2. 回调 (Callback)

我们还可以提供一个回调域，用户也可以为声音的每一个实例设置此回调域。这个回调可以定义成，当声音循环播放或播放了一段时间后才被触发。



```
void SetHandledSoundCallback( Handle_t hSound, CallbackFuncPtr_t pCB );  
void ClearHandledSoundCallback( Handle_t hSound );
```

### 6.1.6 命令 ( Command )

---

“将请求封装成一个对象，从而用不同的请求、队列 (queue) 或日志请求对客户进行参数化，并支持可撤销的操作。”

#### 1. 命令队列 ( Command Queue )

我们可以将一些或全部的音频 API 调用放到一个命令队列中，这个命令队列在每个游戏画面中只被处理一次。通过查看这个队列，可以判断同一个声音在每个画面中是否被多次调用，或者它是否在播放前就被启动和中止了。查看这个队列可以为调试提供重要的帮助。

#### 2. 语音系统 ( Speech Systems )

语音系统得益于命令队列。由于可能有多个人会同时发音，因而可以通过查看请求队列，决定谁是下一个发声对象。有些时候，我们也可以清除请求队列来结束进一步的发音。

```
void PostSpeechRequest( int nSpeakerId, int nPhrase );  
void ClearSpeechQueue();
```

### 6.1.7 备忘录 ( Memento )

---

“在不破坏封装性的前提下，捕获对象的内部状态并使之具体化，以便以后可以恢复对象的该状态。”

#### 暂停和重放 ( Pause and Restart )

当我们要暂停播放所有声音时，可以返回状态的句柄，以后使用这个句柄可以重放这些声音。

```
StateHandle_t PauseAllSounds();  
void RestartAllSounds( StateHandle_t hState );
```

### 6.1.8 观测器 ( Observer )

---

“定义对象间一对多的依赖关系，当一个对象的状态发生改变时，所有依赖于它的对象都得到通知并自动更新。”

#### 动态类型 ( Dynamic Types )

如果我们能将一种类型号传递给正在播放的声音，那么之后，就可以引用这个声音和其他所有拥有相同类型号的声音。例如，如果我们播放某个人的所有拥有相同类型号的声音，

那么之后，可以只通过一次调用就关闭这个人的所有声音。

```
void StartSoundWithType( int nSoundId, int nTypeId );
void StopAllSoundsWithType( int nTypeId );
```

也可以更新所有特定类型的声音。

```
void UpdateSoundsWithType( int nTypeId, const ControlParams_t &cp);
```

也可以将默认优先级和默认音量之类的行为与类型联系在一起。

```
void SetDefaultPriorityForType( int nTypeId, int nDefaultPriority );
void SetDefaultVolumeForType( int nTypeId, float fDefaultVolume );
```

这里假设，当可以使用有限缓存时，存在适当的优先级系统来决定播放哪些声音。

### 6.1.9 大泥球 ( Big Ball of Mud ) ( 也称做 “意大利面条式” 代码 )

我想提的最后一个设计模式是“大泥球”。由 Brian Foote 和 Joseph Yoder 编写的同名论文对模式进行了巧妙深刻的分类，可以自然地融入到项目中。在某种程度上，我们都想避免“大泥球模式”，但相反的是，这又是无法避免的。因此，无论是否懂得模式设计，您都可能执行“大泥球模式”，也就是：编写一次性代码、逐段扩增、保持运转、修修补补，最终弃之不用。

那么，“大泥球模式”如何与音频和游戏编程发生联系的呢？我们的首要目标是避免让音频系统的内部代码成为“大泥球模式”。其次是，避免让余下游戏代码成为“大泥球模式”。在音频 API 中提供一些有用的特性，对代码组织有一些影响。也可以提供 API 功能来帮助调试音频问题。因此，即使陷入了泥潭，音频 API 也可以提供在“泥潭”中导航的功能。

#### 1. 状态函数 ( Status Function )

提供状态函数可以帮助我们找到常见错误。根据特定的请求（如播放的声音号），可以返回特定的值。也可以根据譬如声音状态这种更一般性的请求，返回一个准备显示的字符串。

```
int GetNumberOfSoundsPlaying();
bool IsSoundPlaying( int nSoundId );
string GetAudioStatus();
void GetDescriptionsOfSoundsPlaying( list<string> &StringList );
```

#### 2. 日志函数 ( Logging Function )

我们可以提供一个日志记录系统，用来记录音频系统的行为。这些记录可以定向到不同类型的输出，如标准输出、文件或显示器。我们也可以提供一个细节级别控制，来决定有多少信息需要记录。一个细节级别可能只记录所有音频 API 调用以及传递的参数；另一个细节级别可能记录内部声音缓存分配情况和其他内部状态信息。

```
void EnableLogging();
void DisableLogging();
```

```
void SetLoggingDetailLevel( int nDetailLevel );  
void SetLoggingOutputType( int nOutputType );
```

### 3. 系统禁止 ( System Disable )

有时候判断音频系统是否导致游戏出现故障的最简单方法是完全禁止音频系统。当禁止音频系统后，所有音频 API 调用都不会发挥作用。如果禁止音频系统后问题仍然存在，则可以排除音频系统是导致故障的原因。像这样禁止系统，也可以降低有些游戏调试的复杂性。

```
void AudioSystemEnable();  
void AudioSystemDisable();
```

#### 6.1.10 结论

---

将设计模式当作指导，可以概括出游戏音频 API 的有用特征。无论是在跑道上飙车、在丛林中探险，还是穿越沼泽，我们都希望用户能轻松自如。

#### 6.1.11 参考文献

---

- [Gamma94] Gamma, et al., *Design Patterns*, Addison-Wesley Longman, Inc., 1994.
- [Foote99] Foote, Brian, and Yoder, Joseph, *Big Ball of Mud*, [www.laputan.org/mud/](http://www.laputan.org/mud/)
- [Bilas00] Bilas, Scott, "A Generic Handle-Based Resource Manager," *Game Programming Gems*, Charles River Media, 2000: pp. 68~79.



## 6.2 在采样合成器中声音的同步重用技术

Thomas Engel, Factor 5

thomas.engel@factor5.com

**典**型的合成器只能提供有限的可随时使用的声音（voice），因而需要使用复杂的分配方案充分利用可用声音。只要存在可用的自由声音，就很容易在任意指定时刻分配声音。但如果情况不是这样，那就十分棘手，需要重用声音。这种困难来源于当前激活样本和下一个播放样本之间的幅度差（amplitude difference），这会导致尖锐的爆裂声和卡嗒音出现（如图 6.2.1 所示）。本文将介绍一种技术，可以在不花费大量计算时间和内存的情况下完成声音的重用。

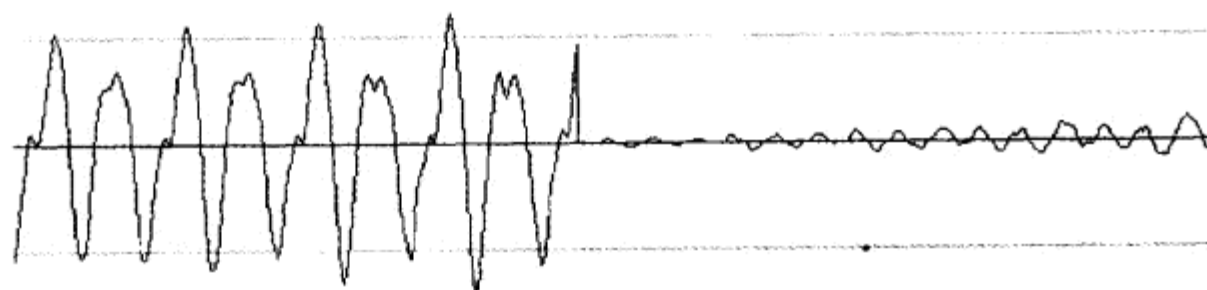


图 6.2.1 样本之间的幅度差将产生令人心烦的卡嗒音

### 6.2.1 存在的问题

如果要在回放另一个声音样本时重用当前激活的声音，可以试着慢慢降低当前样本的音量，并立刻开始回放新样本。这能很好地避免因旧样本突然终止而产生的卡嗒音或爆裂声，并能即时播放新样本。

但不幸的是，这种做法并不十分现实，因为这需要两个声音来正确地处理——并且当这两个声音不再可用时，要同时进行处理。解决这个问题一个可行方案是，为这种情况保留一些声音。但不幸的是，这也有很大的缺点。首先，由于您必须确保一些额外的声音随时是可用的，所以会浪费硬件资源。其次，可能马上重新分配的声音编号将由于此原因而受到额外保留的声音编号的限制。因此，这是一个非常不切实际的解决方案。

您可以通过连续使用如下方法来解决需要两个声音的问题：首先将旧样本渐弱，然后用同样的声音开始播放新样本。这将推迟新声音的开始，而且旧声音音量降低的时间比常规要早很多。这在播放音乐作品时也有问题，因为这种延迟非常显著。不管是否重用声音，延迟所有声音的启动，

会使延迟的效果降到最小。然而，这样的方法会造成音频管道的整体延迟，并且由于增加了复杂性，处理声音变得更困难了。

因此，需要一个更理想的解决方案。

### 6.2.2 解决方案的思路

回答一个简单的问题，可以得到解决方案的基本思路。这个简单问题是：什么是声音（voice）？答案是：声音是气压的变化，也就是由此产生的输出信号幅度的变化。如果输出信号是持续不变的，无论多大的幅度都不能听到。对声音进行重新分配时，导致不期望的幅度突变，刺耳的爆裂声由此产生。即使在不开始新的样本的情况下，停止原来的样本，也能听到这种声音。

现在的思路是，不仅要避免原来样本的任何改变或振动，而且要在输出总和中保持其最终幅度不变。终止单样本的情况如图 6.2.2 所示。

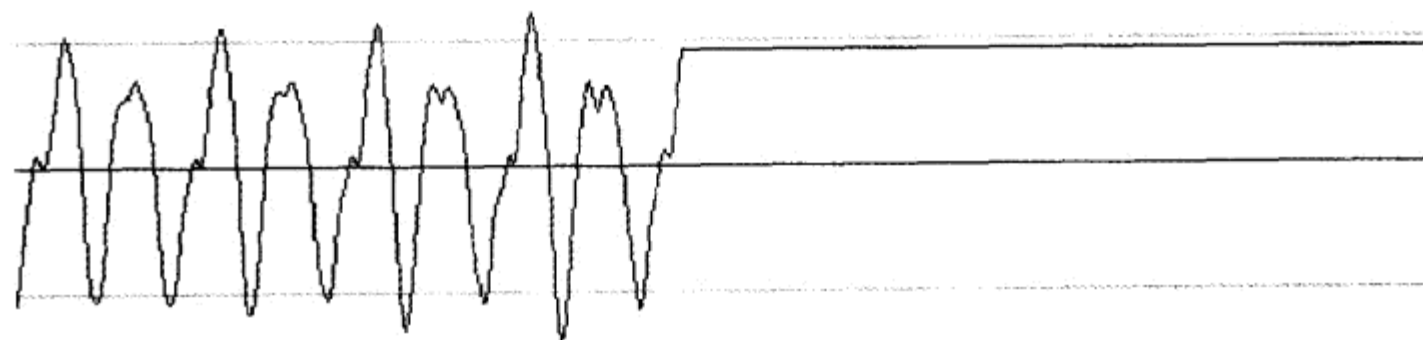


图 6.2.2 保持样品的最后幅度，避免产生不希望出现的爆裂声

这样，由于没有突然的幅度改变，输出就不会产生任何的卡嗒音或爆裂声。

### 6.2.3 解决方案

前面概述的效果可以很好地用来解决我们的难题。中断声音的振动，而不降低它的最终幅度，会使事情简单很多，因为只要开始播放新的样本，那么所有关于获取新样本或处理旧声音音调的逻辑都将被终止。

输出如图 6.2.3 所示。

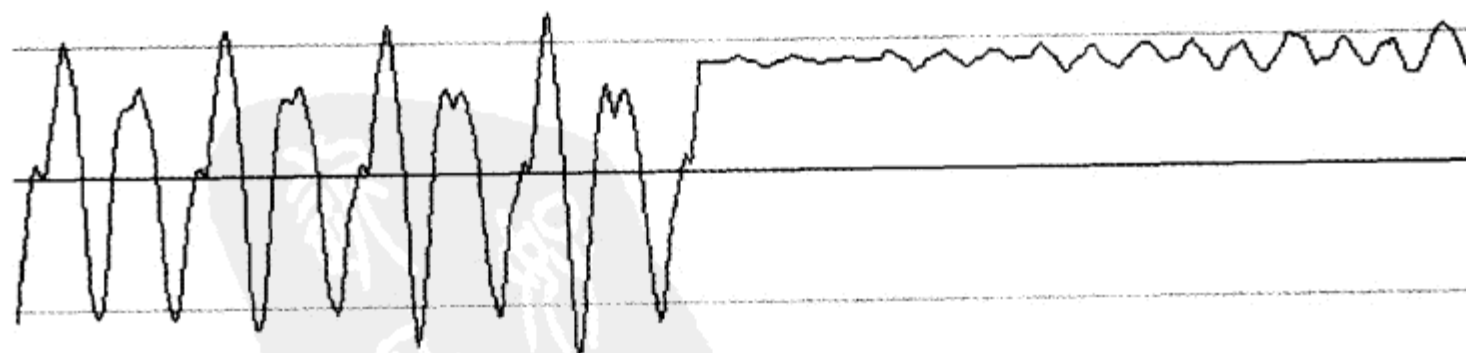


图 6.2.3 以第一个样品的最后幅度开始播放第二个样品

这有个明显的缺陷：简单地中断原来的声音，会给输出信号带来潜在的高 DC 偏移 (DC offset)，这是由于输出样本的最大幅度限制，可能会因为剪辑声音而使输出音质大幅度下降。

有两种简单的方法可以限制（或在大多数情况下完全去除）这样的噪音。首先，可以使用多于 16 位的混合缓存来产生最终输出信号。这就扩大了可内部处理的最大幅度。由于 DC 偏移将和其他能再次降低幅度的声音一起处理，这实际上就明显地减少了噪音 (noise)。然而最重要的步骤是，快速地将 DC 偏移降低至零；主要任务是，必须让终止的旧样本渐弱（如图 6.2.4 所示）。

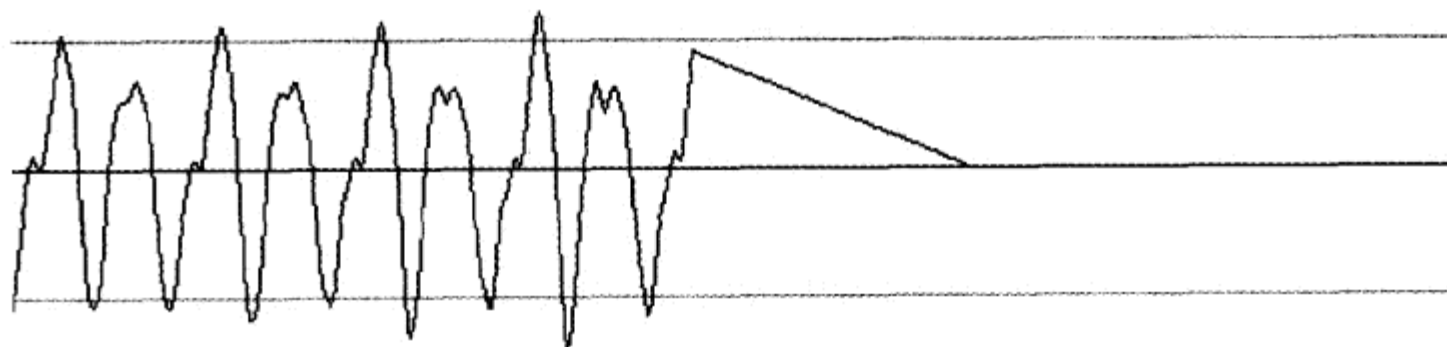


图 6.2.4 将最初的样品渐弱至零

这仍然要求对每个声音都要有单独的逻辑来处理终止的旧样本，避免声音数量的限制，以可以立即重用它们。事实果真如此吗？实际上，通过一个非常简单的处理器，可以处理重用声音引入的所有 DC 偏移。

逻辑过程如下：每次重用声音时，由特定声音上的旧样本混合到总和中的最终值将被添加到总 DC 偏移值上。一旦所有的声音都被处理后，这个值就会反过来添加到输出信号中。

要去除前面讲述的这种必然但却不希望的 DC 偏移，那么，每过一段时间，就以每个样本为基础，降低 DC 偏移总和。这种做法很有效，并且只要不是操作过快，就不会产生任何更多的噪音。5 ms 的时间足够将最大振幅的 DC 偏移降低到零（如图 6.2.5 所示）。

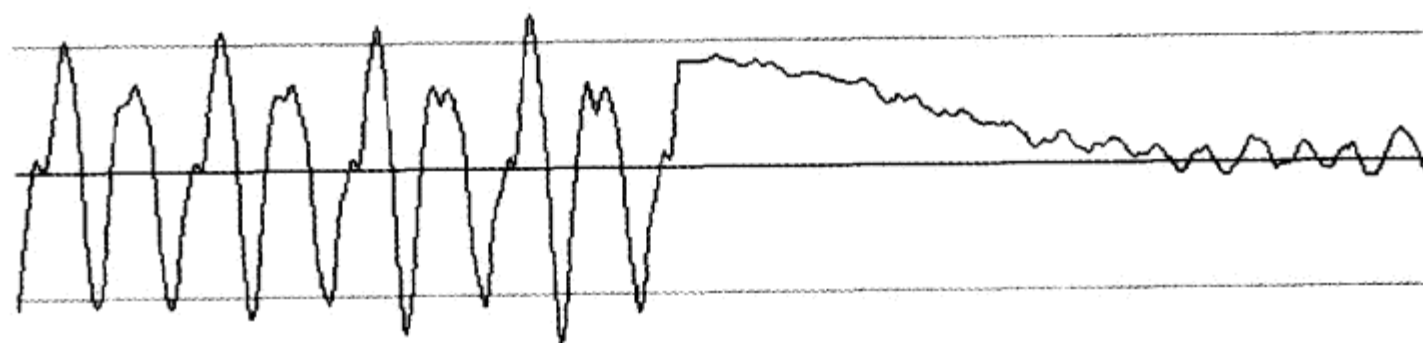


图 6.2.5 在播放第二个样品时，将最初的样品渐弱至零

## 6.2.4 结论

实际上，理论上 DC 偏移的负面效应几乎是觉察不到的。只有在用至少 24 位的混合缓存，重用大量高幅度的声音，才会达到一些可觉察的效应。

同时，这种方法的益处也是巨大的。我们可以在没有任何延迟的情况下，立即对声音进行重用，输出信号中也不会有任何可听到的噪音。使用这种方法的代价是最小的。添加一些



额外操作后，这种方法可以用在具备多声道硬件样本回放功能的系统中。按照本文讲述的方法，只需要留出一个声音来容纳加在由 CPU 产生的样本数据流中的直流偏移值就可以了。



## 6.3 软件 DSP 效果

Ian Lewis, Acclaim Studios  
ilewis@acclaim.com

这篇技术精粹将介绍 DSP（数字信号处理）的基础概念和技术，这些技术包括滤波（filtering）、卷积（convolution）、延迟和插值（interpolation）。可以毫不费力地将这些技术应用到现在的 PC 机和下一代控制台上（参见表 6.3.1）。

表 6.3.1 DSP 技术	
DSP 技术	用 途
滤波	封闭和 3D 声效
	物理建模效果，如引擎的拉动
卷积	头部相关传输函数(Head-Related Transfer, HRTF)
	精确地模拟真实声音空间
延迟	墙壁或地面传来的回音
插值	音调偏移
	采样速率转换

### 6.3.1 滤波

滤波实际上是卷积的一种特殊情况，但由于它比一般的卷积更容易快速实现，所以独成一体。

滤波音频信号就是加强或削弱某些频率成分，就像是立体声系统上“低音”和“高音”的旋钮一样。这种影响在波形图上很难看出，但耳朵却能立即感觉到。例如，使用低通滤波器会使声音听起来低沉，而使用高通滤波器则使声音变得单薄。

滤波的一般形式是一个循环过程，在此循环中，输入样本乘以一组数值（滤波系数），再进行累加。第一个滤波系数与当前样本相乘，第二个和上一个样本相乘，依此类推，代码如下：

```
for (k = 0; k < nCoefficients; k++)
{
    output+=coefficient[k] * (*(input-k));
}
```

这种滤波器叫 FIR 或有限脉冲响应滤波器（Finite Impulse Response filter），因为在输入值衰减到零后的一段可知时间内，它的输出值总会衰减

到零。一种在声学上略欠稳定，但却更强大的滤波器是 IIR 或无限脉冲响应滤波器（Infinite Impulse Response filter）。这种设计将滤波器的输出反馈给输入，形式如下：

```
for (k = 0; k < nCoefficients; k++)
{
    output += coefficient[k] * input[n-k] +
        feedback_coefficient[k] * previousOutput[n-k];
}
```

当然，在计算完每个样本后，previousOutput[]数组的内容会上移 1，当前输出被放到数组的最后。这看起来更加复杂，但好处是系数的数量比等价的 FIR 实现小得多。

尽管现在有几百种滤波器，但大多数都是面向音乐应用的，比游戏中需要的更为复杂。示例代码使用了音乐 DSP 源代码文档中的简单钱伯伦滤波器（Chamberlin filter）[DeJong01]。在它的完全形式中，提供了高通滤波器、低通滤波器、带阻滤波器（notch filter）和带通滤波器（bandpass filter），这种滤波器开销非常小并且容易实现。

滤波器可以很好地和 SIMD 硬件，甚至是和矢量硬件一起使用，因为它们就是乘法累积串。

### 6.3.2 卷积（convolution）

卷积是将获得的两个信号相乘，使第一个信号带有第二个信号的特性的过程。例如，假如说在一座旧大教堂里装有水平集（level set），可以在附近任何一座歌特式大教堂里真实地将“脉冲”（一种很短的声音尖峰信号，就像气球的爆炸声）录制下来，混合这种脉冲与游戏中的声音——游戏的声音听起来就像真的从那座大教堂中发出。不幸的是，求卷积的开销相当大，主要是因为每个输出样本都是输入样本与所有脉冲样本相乘的总和。要理解这些代码，可以将脉冲文件中的每个样本看作前述 FIR 滤波器中的一个系数。这只真正适用于很短的脉冲。

MIT 媒体实验室中的头部相关传输函数（HRTF）集就是一个理想的短脉冲集[Gardner94]。使用合适的 HRTF 的卷积，可以产生一种置声音于三维空间的错觉。

一种比较快速的卷积方法需要使用快速傅立叶转换（Fast Fourier transform）将两种信号都转换到频域（frequency domain）内（就像在分光镜下，每个样本表示一个频率，而不像是示波镜里的幅度）[Press92]。如果已经将样本存储在频域内了，使用这种方法会比较有效。有些压缩方法的确使用了 FFT 或它的某一变种，例如，MP3 压缩使用的就是一种叫做离散余弦转换（Discrete Cosine Transform, DCT）的变种。

### 6.3.3 延迟

数字延迟可能是最容易产生的 DSP 效果。它只用很少的 CPU 时间，但需要一些内存。它的算法很简单：将每个样本的一份拷贝放到一个队列（queue）或循环缓冲（circular buffer）中，然后将最早的样本从缓冲取出，并与输出混合在一起。

```
DelayBuffer.push(input);  
Output=input + DelayBuffer.pop() * delayGain;
```

缓冲的长度决定了延迟的长度。通常，您会希望依某个信数（如上述例子中的 `delayGain`）放大延迟音量。例如，一个简单的延迟可以为室外大场景增添回声。

更有趣的延迟是再生延迟（regenerative delay）。在再生延迟中，放大的延迟输出样本被反馈回延迟缓冲，产生一个多重、能更准确模拟自然空间的回音效果。由于每个样本的增益在每次循环时不断减少，最终，回音将消失。延迟消失的时间由每个样本中使用的衰减量决定。衰减量多，会有短暂的击掌回音效果，仿佛站在峡谷边或大建筑物墙边一样；衰减量少，则得到一阵密集的回音，更像在教堂或回音洞中听到的回音一样。

更复杂的混音效果是多插座混响（multi-tap reverb）。在多插座算法中，输入样本的拷贝不仅被加到延迟缓冲的尾部，还被插入到缓冲的其他位置（这就是所谓的“插座”）。这可以得到更密集的回音和更真实的效果。例如，在一个矩形大房间中录制的声音可以在延迟缓冲中插入 4 次，一次插入针对一面声音反射墙。缓存中每个“插座”的准确位置，根据声音源离墙的距离来动态决定。

空间感的基础来源于 DSP 混响算法，这是基础延迟理论更复杂的变异。一般来讲，由于混响算法开销太大，所以不能在软件中实现。由于在 PC 声卡（如创新公司的 EAX 扩展卡）和下一代游戏控制台上实现硬件混响处理器的势头不断增长，本文不对混响算法进行深入地描述。

#### 6.3.4 插值（interpolation）

插值是一门根据已知样本值，确定不存在样本值的技术。这看起来有点深奥，但其原理在游戏编程的很多领域经常使用的。例如：双线性纹理滤波（bilinear texture filtering）就是插值的一种方式，其中可以只使用 4 或 5 个质素（texel）来覆盖屏幕上 10 或 20 个像素。重新采样是插值在音频中主要的用途。压缩音频的最有效方法是减少采样率，特别是当音频有少量高频率成份的时候。

采样率转换（和音调偏移大体上相同）可以使用几种算法。下面按开销从少到多的顺序列出：

- **成倍采样** 适用于目标采样率和源采样率之间的比率是 2 的倍数时。
- **平均采样** 适用于计算的目标样本为两个最近点数值的平均数时。
- **线性插值** 与平均采样相似，不同的是，平均数是经过加权的——点越近，加权越多。
- **三次插值** 取最近的 3~7 个点（据报道，取 5 个点的效果最好[DeJong01]），使用样条函数逼近，目标点的值等于那一点样条的值。

三次插值确实比重新采样的其他形式效果好，但令人吃惊的是，所有其他形式的效果相同。当重新采样率不是 2 的幂的情况下，如果 CPU 的运行速度有限，那么简单的平均采样方法是最好的选择。对于源采样率与目标采样率的比值介于 0.75~1.5 之间的音调偏移，将获得合理的效果。

三次采样可以得到较高的比率，而不会产生难听的声音，但简单重新采样算法不可能避

免可怕的“米老鼠”声效。要避免“米老鼠”声效，则要进行共振峰修正音调偏移，这涉及到一些非常复杂的计算，这已超出了本文的讨论范围。

### 6.3.5 参考文献

---

[Bores01] “Bores Signal Processing. Introduction to DSP,” [www.bores.com/courses/intro/](http://www.bores.com/courses/intro/) (4 March 2001).

[DeJong01] DeJong, Bram, “The Music-DSP Source Code Archive,” [www.smartelectronix.com/musicdsp/](http://www.smartelectronix.com/musicdsp/) (4 March 2001).

[Gardner94] Gardner, Bill, and Martin, Keith (1994). “HRTF Measurements of a KEMAR Dummy-Head Microphone.” *MIT Media Lab*, <http://sound.media.mit.edu/KEMAR.html> (4 March 2001).

[Iowegian99] Iowegian International Corp. (1999), DSPGuru, [www.dspguru.com/](http://www.dspguru.com/) (4 March 2001).

[Kosbar98] Kosbar, Kurt L. (1998), “Introduction to Digital Signal Processing (DSP),” [www.siglab.ece.umn.edu/ee301/dsp/intro/index.html](http://www.siglab.ece.umn.edu/ee301/dsp/intro/index.html) (4 March, 2001).

[Press92] Press, William H., et al, *Numerical Recipes In C*, Cambridge University Press, 1992, [www.ulib.org/webRoot/Books/Numerical\\_Recipes/](http://www.ulib.org/webRoot/Books/Numerical_Recipes/) (4 March 2001).

[Sprenger01] Sprenger, Stephan M, “The DSP Dimension,” [www.dspdimension.com/](http://www.dspdimension.com/) (4 March 2001).



## 6.4 数字音频的交互式处理管线

Keith Weiner, DiamondWare Ltd.

keith@dw.com

### 6.4.1 简介

数字信号处理 (Digital Signal Processing, DSP) 的算法很久以前就已经出现。但游戏却没有充分地利用它们, 大概是因为它的高 CPU 开销。不过, 现在吉赫级的处理器相当便宜, 而且越来越便宜。这使得一定数量的音频处理成为可能; 甚至在图形引擎、AI 引擎和物理引擎共享 CPU 时间后, 仍有处理音频余地。

游戏 (以及其他应用) 具有音频处理管线的需求。目标是为某种声音指定一套带有相应参数的 DSP 函数。其他声音则可能通过一套不同的 DSP 函数, 或者使用不同的参数集。

图 6.4.1 为音频系统的框图。完整的结构将从任何源流出声音, 在必要时解码、解压、处理、混合, 然后输出。本篇文章描述了标记为 “Processing” 的结构块。

本篇文章介绍了一种轻重量级, 但通用的结构。“通用” 在这里指两个方面。首先, 它允许任何类型的处理, 甚至是速度更改器。速度更改器有一定难度, 因为必须改变缓冲长度。如果希望任意更改参数, 则难度更大。其次, 它完全与 OS 独立, 而且在 Linux 和 Windows 上可以轻松移植使用。这篇文章演示的结构已经应用于设备驱动程序、通用音频 API/引擎以及电话媒体控制组件中。

#### 限制

在构思并设计一个系统时, 抓住重点的一个好方法是定义问题以及解决方案的所有限制。下面将所有需求列出, 然后讨论每一项。

- 支持  $N$  个声音。
- 每个声音可以经过  $N$  个 DSP 函数。
- DSP 函数的任何参数在任何时间可以改变。
- 任何函数可以强制改变长度。
- 高效。这是针对 PC 上的游戏, 而不是超级计算机上的科学研究。
- 使之更易于编程 DSP 函数。



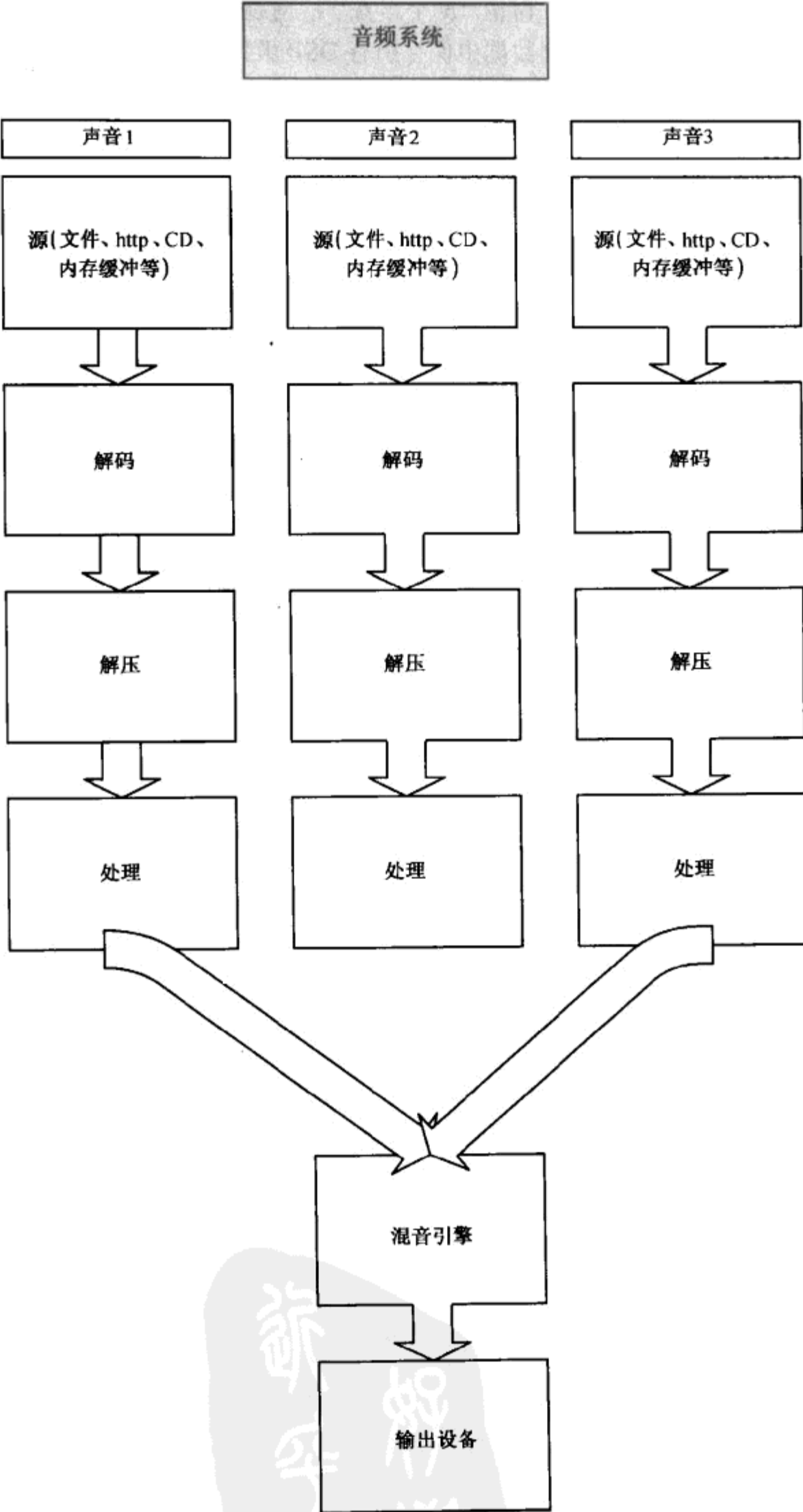


图 6.4.1 音频处理系统

下面详细讨论每一项。首先，讨论“ $N$ 个声音”的问题。单声引擎在1994年以后就被淘汰了。如果我们在每个声音的实例数据中保存所有DSP函数和参数信息，则可以处理无限的声音，而不会带来额外的复杂性。

下一个要求是每个声音支持 $N$ 个DSP函数。这是本篇文章的核心。如果希望让一个声音只经过单个DSP函数，则不需要管线。

参数可以在任何时候改变。这是交互式的核心。在游戏中，声音设计师不知道将发生什么，也不知道什么时候发生。在任何时间，用户可以进入具体管道，然后某定声音必须进行混响。

也许您不会太频繁地降速或加速正播放的声音，但有时会这样做。只要我们是构建一个通用的引擎，最好让它支持这种运算。在管线上下文中，速度改变DSP函数是输出的样本数与输入的不同的函数。放慢音乐的DSP函数例子只是速度更改器的一种。如果混响函数生成的输出可以穿过输入的末尾（换句话说，它允许最后的混响在停止前渐弱到静音），则该混响函数可以归类到速度更改器中（虽然它实际上更应算是一种长度更改器）。

音频系统必须高效。可以认为，管线本身不会消耗大量CPU周期，但要小心不要进行不必要的内存复制。例如，如果DSP函数返回“无工作可做”，也不要执行内存复制。

任何通用API的目标是保存工作，而不是创造工作。音频引擎只构建一次。DSP函数有多种类型。我们应该完成尽可能多的辅助工作，让每个DSP函数实现一个算法，而不是实现很多算法。

## 6.4.2 讨论

在这一节中，我将针对一个声音的情况进行讨论。直到混合阶段，才会涉及到多个声音的支持。

再次考察图6.4.1，可知我们得到的音频输入来自于解压器（若不存在解压器，则来自于输入流管理器）。输出进入混合引擎（它可能为DirectSound，或其他系统）。

这是不对称的。虽然输入端相对灵活（它应该能够为我们提供所需要的音频样本），输出端不灵活。混音引擎需要取自流的一定数量样本。它不能再处理任何更多的样本，而且如果处理少于此数量的样本而继续混音，则没有声音，即存在断裂。

一条必须重述的重要规则就是，音频处理管线输出必须由混音引擎的需要决定，它决定输出多少样本。

输入流只有一种类型的不灵活：EOF。当流到达末尾时，处理引擎需要知道。

因此，混音是以固定的16384个样本的缓冲来完成的。这意味着它要一直等到需要这么多的样本为止，然后触发音频处理。所以，音频处理管线必须输出16384个样本。

对于输入，管线需要多少样本？我们不得而知。到日前为止，我们只知道链中的最后函数必须输出16384个样本。那么，函数需要多少输入的样本呢？

我们不知道，但可以找出答案。可以查询得知。混音器希望16384个样本。链中的最后函数，通过它的当前参数集，提出它希望24576个样本的输出来生成16384个样本的输出。为了决定从流中获得多少个输入样本，继续向后遍历DSP函数链。向后遍历是因为知道最终输入样本数，但不知道起始的输入样本数。

一旦完成了此阶段任务，就可以从输入流获得所要求的样本数，让它流经链中的第一个函数。可以取出它的输出，并流经第二个函数，如此继续。从理论上讲，最终会得到发送到混音引擎的正确样本数。

如果假设不成立呢？如果某个 DSP 函数生成的样本数少于应该生成的样本数呢？链中的下一个函数必须赋予新的、更小的缓冲。它必须生成尽可能小的输出。最终可能存在不足量。将这个不足量作为第二次迭代的期望输出，沿链向后查询，然后第二次处理数据。继续进行，直到用完输入样本，或者充满了输出缓冲。基于效率方面的考虑，DSP 函数的实现方式应该是让它们不轻易误报所期望的输入要求。只有计算准确要求的输入样本数需要实质性的处理时，这个功能才有用（例如，当 DSP 函数必须实际处理音频，才能准确知道它需要多少样本时）。

如果 DSP 函数“吃掉”了比所需少的样本，情况该如何呢？这是编程中的比较处理。管线必须在链中的每对 DSP 函数之间提供一个缓冲。函数#2 和函数#3 之间的缓冲容纳函数#3 中的下溢数据。

如果 DSP 函数生成了太多的输出样本，又该怎么办呢？事实上不应该这样做。认真地讲，处理这种情况的有效办法是缓冲它，不过，本篇文章没有相应的演示代码。

这就出现一个问题：DSP 函数不可交换，即顺序的问题。这不仅是我刚才所提的缓冲系统的限制，而且也是处理过程的本质所在。如果先更改声音的音调，再进行混响，其结果与先混响再更改音调不同。因此，不必担心我们提出的结构会同时强加此要求。

### 细节

想像有一个音频处理管线，只要混音引擎需要更多的数据用于声音通道时就会调用它。如果混音引擎由中断驱动，可以实现为一个回调，或者在混音缓冲被查询时作为一个常规函数。

处理引擎必须调用输入流管理器从声音通道请求原始数据。它可能返回一个全缓冲，或者声音到达末尾时较小的缓冲。

处理引擎可以通过若干信息调用每个 DSP 函数。到目前为止，我们讨论了决定生成目标输出所需要的输入样本数的查询，当然，也讨论了输出调用本身。

还要添加两个：firstcall 和 lastcall。Firstcall 为发送到 DSP 函数的第一条信息，让它分配实例数据；而 lastcall 允许它释放数据。

要实现一个 DSP 函数插入到这个结构中，程序员必须提供单个入口点来正确处理 firstcall（初始化）、lastcall（关闭）、查询以及输出信息。赋予这个入口点一个结构，此结构包含命令数、输入和输出缓冲和长度以及其他两个项：其中一个为特定 DSP 参数的指针，它从应用向下传递（特定 DSP 参数的例子有混响中的房间大小、速度更改器中的比率等等）；另一个项为 DSP 函数本身分配和维持的私有实例数据。

实现这个引擎在一定程度上更复杂，但这是我们的目的。让这一块智能化后，许多 DSP 块就能快速、轻易地实现。

处理引擎使用 3 个缓冲：一个源缓冲、一个目标缓冲以及一个累加缓冲。在处理过程中，缓冲交换角色，因此引擎使用 buf[0]、buf[1] 和 buf[2]。原因是 DSP 函数#1 的目标（输出）变成 DSP 函数#2 的源（输入）。

与之相似,如果链输出中的最后 DSP 函数小于混音引擎所需要的数量,它的目标缓冲变成累加缓冲。在下一次链迭代后,累加缓冲预先挂起为新目标缓冲。

### 6.4.3 代码

下面是主要入口点的源码。

```
void dsp_ProcessAudio(dsp_AUDIO *audio)
{
    audio->src=0;
    audio->dst=1;
    audio->acc=2;

    bfbuf_WIPE(audio->buf[audio->src]);
    bfbuf_WIPE(audio->buf[audio->dst]);
    bfbuf_WIPE(audio->buf[audio->acc]);

    audio->samps = 0;

    while ((audio->samps < audio->sampsneeded) && (!audio->done))
    {
        SWAP(audio->acc, audio->dst);

        QueryDSPIn(audio);
        stream_GetData(audio);
        DoDSP(audio);

        bfbuf_MoveAll(audio->buf[audio->dst], audio->buf[audio->acc]);
    }
}
```



在使用以上代码前,还有一点需要解释。bfbuf 是我创造的一个数据结构。回填缓冲 (backfill buffer) 是设计用来让预先挂起数据变得容易的内存缓冲。其内容总是在结束时进行调校刷新,所以用一个简单的 `memcpy()` 来完成预先挂起。BFBUF.C 的源码请浏览附带光盘。

在处理 DSP 链本身的下溢缓冲时,预先挂起是我们所必需的。因此,我实现了外部循环,它利用预先挂起从前一次迭代中获取缓冲部分,并将其添加到当前缓冲部分中。正如在先前部分缓冲上附加当前部分缓冲一样简单。当您有了一把新“锤子”,而且感觉很精巧,您就会倾向于将任何东西看作一颗“钉”。



让我们考察这个函数。它的惟一参数是结构的指针 (结构的定义请见光盘上的代码)。有 3 种回填缓冲,它们初始化为空,并赋予初始语义。

我们使用了一个循环,它在用完源数据或输出缓冲充满前一直迭代。对于第一次迭代,交换无意义,它什么作用也没有。

剩下问题很简单。首先，决定需要多少输入样本，获得这些样本，并进行处理。然后，预先挂起前一迭代中的累加样本（如果有的话）到当前目标上。

`stream_GetData()`是对输入流管理器（DSP.C 外部）的调用。这个函数必须将样本放在目标缓冲中，理解这一点很重要。我们将看到，`DoDSP()`在开始时交换源与目标缓冲。

下面介绍查询函数。

```
static void QueryDSPIn(dsp_AUDIO *audio)
{
    dsp_PARAMS params;
    dsp_DSP *dsp;
    DWORD underlen;
    DWORD x;

    params.cmd = dsp_QUERYACTIN;
    params.actualin = audio->sampsneeded - audio->samps;

    for (x=0;x<audio->numdsps;x++)
    {
        dsp = &(audio->dsp[audio->numdsps - x - 1]);
        params.dstlen = params.actualin;

        (*dsp->callback)(&params);
        underlen = bfbuf_GETDATALEN(dsp->underbuf);

        dsp->inreq = (underlen < params.actualin) ?
                    params.actualin - underlen : 0;
        dsp->outreq = params.dstlen;
    }
}
```

注意函数开始时是如何知道剩下多少样本的（`audio->sampsneeded - audio->samps`）。让循环中的代码从前一次迭代中找到它。也可以用一个 `goto` 语句在中途进入循环进行编码，但这种编码方式似乎不是雅致的解决方案。

循环本身以逆序遍历声音使用的所有 DSP 过滤器。对于每一个 DSP 过滤器，它设置一个参数来指定期望的输出样本数，并进行调用。然后，它减去前一次迭代处理中保存在下溢缓冲中的样本数。如果前一次迭代留下足够的样本，结果可能为 0。在这里，过滤不再需要输入。它所需要的所有输入已经在下溢缓冲中准备就绪！

每个 DSP 保存它需要多少输入样本，以及期望的输出（链中的下一项要求的样本数）。

本篇技术文章的精粹是 `DoDSP()` 函数。

```
static void DoDSP(dsp_AUDIO *audio)
{
    bfbuf_BUFFER *srcbuf;
    bfbuf_BUFFER *dstbuf;
    dsp_PARAMS params;
    dsp_DSP *dsp;
    DWORD x;

    params.cmd = dsp_OUTPUT;
```

```

for (x=0;x<audio->numdsps;x++)
{
    dsp = &(audio->dsp[x]);

    SWAP(audio->src, audio->dst);

    srcbuf = audio->buf[audio->src];
    dstbuf = audio->buf[audio->dst];

    bfbuf_MoveAll(srcbuf, dsp->underbuf);
    bfbuf_SetDataLen(dstbuf, dsp->outreq);
    bfbuf_GetBufStart(srcbuf, &(params.lsrc), &(params.rsrc));
    bfbuf_GetBufStart(dstbuf, &(params.ldst), &(params.rdst));

    params.srclen      = bfbuf_GETDATALEN(srcbuf);
    params.dstlen      = bfbuf_GETDATALEN(dstbuf);
    params.dspspecific = dsp->dspspecific;
    params.privdata    = dsp->privdata;

    (*dsp->callback)(&params);

    if (params.nowork)
    {
        bfbuf_WIPE(dstbuf);
        SWAP(audio->src, audio->dst);

        params.actualout = params.srclen;
    }
    else
    {
        if (params.actualout < bfbuf_GETDATALEN(dstbuf))
        {
            bfbuf_ChokeUp(dstbuf, params.actualout);
        }
        if (params.actualin < bfbuf_GETDATALEN(srcbuf))
        {
            bfbuf_Eat(srcbuf, params.actualin);
            bfbuf_MoveAll(dsp->underbuf, srcbuf);
        }
        else
        {
            bfbuf_WIPE(srcbuf);
        }
    }
    audio->samps += params.actualout;
}

```

源与目标缓冲的初始交换看起来很有趣。每个 DSP 函数的输入（源）是前一个函数的输出（目标）。这也使处理 DSP 执行无工作可做的情况变得容易。

这里还有一些变量设置，然后 4 次调用这个神秘的回填缓冲（backfill buffer）管理器。



首先，从前一次迭代取得所有下溢样本。注意，它按逐 DSP 来保存。DSP 函数#1 的下溢不同于 DSP 函数#2 的下溢，而且必须分离。

接下来，设置目标回填缓冲的期望大小。如果它是错误的，则在以后修正它。最后，获得当前缓冲起始处的指针。请注意，分开保存左右指针。这使得 DSP 处理更容易，而且最终移到多通道音频中（例如 5.1 节）也更容易。

我们进行一些直接的参数设置。注意如何为 DSP 提供特定函数参数的指针以及声音的实例数据。

调用后，如果函数不工作，则设置目标缓冲为空，并交换目标与源缓冲（因此目标现在为其中具有样本的缓冲），以此来准备下一次迭代，它将再次进行交换。

否则，函数将执行工作，因此，我们要检查 DSP 生成的缓冲是否小于满缓冲。Choke Up 是惟一必须进行内存移动的 BFBUF 运算。

然后，检查函数的消耗数据是否小于提供给它的所有输入数据。如果是这样，则作为下溢缓冲。否则，只将源缓冲标记为空。

#### 6.4.4 额外注释

这不是最简单的代码，但却雅致并有效，这正我们所追求的。没有周期浪费；除非 DSP 生成的数据太少，而这也只用再次迭代整个链以得到足够的数据外。

注意每次 `dsp_ProcessAudio()` 是如何被调用的，所有 DSP 函数的参数可能发生变化。这些代码形成了音频处理引擎的核心，它符合我们在前面列举的所有标准。

构建它不仅可能，而且并不像看起来的那么难。



正文中忽略了两个函数（收录在附带光盘中）。在系统中添加新声音时，必须调用 `dsp_NewSound()`。当从系统删除声音时，必须调用 `dsp_DeleteSound()`。这一点很重要，因为这些显示的代码不太明显。存在一个对流管理器的调用，当流到达末尾时，调用设置 `audio->done`。

不管是什么代码调用，`dsp_ProcessAudio()` 必须检查这个位，然后在适当时调用 `dsp_DeleteSound()`。技巧是要一直等到声音完成（即 `audio->samps` 为 0）。DSP 函数（如回音）可以保持它们自己的内部缓冲。不管它们如何决定何时“吃掉”输入样本、不生成任何样本（当查询输出以及提供样本比需要的少时），它们都必须尽其所能输出样本。按这种方式，DSP 引擎知道何时流到达末尾、DSP 函数链产生无输出、完成声音回音的混响以及完全删除声音，而不导致可觉察的截短。

从结构的立场看，这诚然有些脆弱。但对于 DSP.C 内的这个问题，可以设计一个更雅致、独立的解决方案。一个相关的限制是，`dsp_AUDIO` 结构以及它的 `dsp_DSP` 数组同时在外部分配。在 DSP 引擎之外，与 DSP 相关的内存结构、参数等等太多。这样做主要是要保持代码的整洁，而且更易于理解。

有两个函数上的限制要讨论（前一个被认为是结构上的）。首先，当 DSP 函数溢出时，此代码不能处理这种情况。原因是没有函数输出多于要求的数据。也可以认为，那些新函数只可以维持一个内部缓冲。这并不能令人满意，因为我们已经竭尽全力确保 DSP 函数为小型、

轻量级而且易于使用的。真正的答案是，它将大大增加代码规模和复杂性（包括 BFBUF.C），并提供一个（相对的）临界值。这段代码不是真正的健壮商业实现。例如，在放入 DEBUG 错误检查之前，添加溢出保护这样的代码（包括决定多大才是足够大等等）并无意义。

另一个不具备的功能是，创建声音后无法添加或删除 DSP 函数。只要所有的 DSP 函数有一个参数来指示“passthru”，就可以解决这个问题。所以，只需要确保通过处理于生存期、可能想要使用的所有 DSP 函数来创建每个声音。这并非无效，可以正常发挥作用（尽管有些笨拙）。

这可能是不太理想的杂合，但只要您为它们提供了 firstcall 处理的方式，就可以实际添加新的 DSP 函数。

#### 6.4.5 结论

这篇文章的实质是演示如何构建基本的 DSP 管线，支持流经声音通道，通过 DSP 过滤链，每次过滤可以对它需要的流进行处理。

为了支持那些希望在处理时更改缓冲大小的 DSP，我们需要一个查询调用，询问 DSP“是否我们需要  $N$  个输出样本，估计需要多少输入样本？”按逆序进行该查询。



另外，添加了对两个下溢条件的支持：生成太少的数据以及吃掉太多的输入。在新数据结构回填缓冲的帮助下，它们不难实现。概念最难掌握，但其代码简单易懂（相关代码可以浏览附带光盘）。

通过这个引擎以及一些胶合逻辑，可以进行游戏所需的任何类型的音频处理。实际上，它可以支持一个完整的专业音频工作室，策划音乐作品和表演。



## 6.5 游戏中的基本音乐音序器

Scott Patterson

scottp@tonebyte.com

本文描述了如何生成基本音乐音序器 (sequencer)。首先, 从比较音乐流和音序开始。然后, 介绍 MIDI 音乐语言中的一些重要概念, 概括基本计算机音乐语言实现, 然后讨论与定时和合成控制相关的编码问题。

### 6.5.1 音乐流与音序

游戏中的音乐一般按两种方式之一进行回放: 作为样本数据的流或者作为音频合成系统的指令序列。音乐回放也可能是这两种方法的综合。这些方法各有优缺点, 它们依赖于所使用的硬件系统。为了综合起来, 还要解决合成问题。

#### 1. 流方法的优点

- 音频硬件需求极少。
- 音乐编制与游戏集成容易。
- 音乐质量仅局限于采样率以及流数据的通道数。

#### 2. 流方法的缺点

- 每次按相同方式插入音乐数据。
- 流起始和终止之间存在延迟时间。
- 当循环一个流或切换流时存在延迟问题。
- 存在大内存缓冲需求。
- 在回放期间可能冻结了一定资源, 如 DMA 硬件和存储设备硬件。
- 如果流数据没有压缩, 它将占用大量存储空间。
- 如果流数据经过了压缩, 它将占用硬件或软件解压资源。
- 根据硬件的限制, 也许不能同时渐弱或渐强两个不同的流, 因此, 惟一的选择可能是渐弱到无声, 或者从无声渐强。

#### 3. 音序方法的优点

- 起始或终止音乐时不存在延迟问题。
- 循环或切换音乐时不存在延迟问题。
- 如果音频内存可用, 则样本数据内存不会与游戏数据内存竞争。

- 回放音乐时可能不需要 DMA 硬件和存储设备这些硬件资源，而且在任何时候可以释放这些资源以供游戏使用。
- 根据相当小的样本数据集可以生成大量音乐风格和变化。
- 音乐数据可以在运行时动态改变，创建独特的游戏和音乐交互。

4. 音序方法的缺点

- 音乐中的声音可能与其他游戏音频所需要的声音竞争。
- 音乐编制以及与游戏的集成更为复杂。
- 音乐质量可能受到合成 (composite) 功能的限制。
- 音乐质量可能受到可用样本内存的限制。

6.5.2 核心计算机音乐概念

下面介绍一些核心计算机音乐概念，同时构建我们的音乐命令语言。

1. 事件块

音乐可以描述为一段时间内的事件 (event) 或命令发生。因此，可以通过“事件块”构建音乐命令语言。事件块由 3 个元素组成：时间、事件类型以及事件细节。我们将保存相对于一个事件相关的时间，称之为增量时间 (delta-time)。事件类型用一个数字来标识，事件细节为 0 或更多的参数，它们由事件类型定义 (如图 6.5.1 所示)。

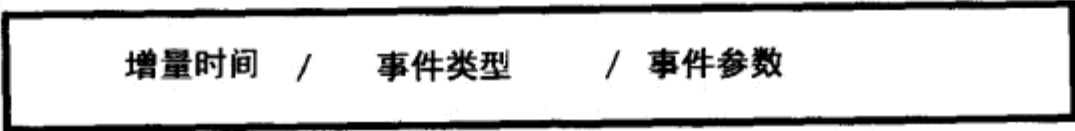


图 6.5.1 事件块

使用事件块的序列 (sequence)，可以描述任何事件和时程 (time interval) 的组合。所以，构建计算机音乐语言的任务就是选择如何保存增量时间，选择支持的事件类型以及选择提供给每种事件类型的参数。

2. 乐器数字接口 (Musical Instrument Digital Interface, MIDI) 音乐

MIDI 规范已经出台很久了。大部分音乐合成和定序软件提供了通过 MIDI 输入和输出端口进行录制和回放的功能。这些软件产品也提供了保存和载入 MIDI 文件的选项。因为我们要创建自己的计算机音乐语言，所以最好从学习 MIDI 规范及其文件格式开始。然后将 MIDI 文件数据转换成自定义的数据。

“接口”一词在 MIDI 中指乐器间通信的规范。在制作游戏的音乐音序器 (sequencer) 中，并不像在特殊的机器上插入音乐一样对乐器间的通信感兴趣。因此，这里对 MIDI 的介绍只包括与游戏中音乐音序器直接相关的部分。

3. MIDI 事件类型

MIDI 规范定义了若干种事件(event)，称做通道声音消息(channel voice message)。MIDI 最先作为一种虚拟钢琴键盘通信语言，因此，起音力度(attack velocity)、音调轮(pitch wheel)以及键压(key pressure)都指应用于虚拟钢琴键盘上的动作(参见表 6.5.1)。

表 6.5.1 MIDI 通道声音消息

事件类型	事件参数
音符关闭	音符号 释放力度
音符打开	音符号 起音力度
音调轮	音弯 LSB 音弯 MSB
控制改变	控制器 ID 控制器值
程序改变	程序号
键压	音符号 压力值
通道压力	压力值

通道声音消息是 MIDI 定义的核心音乐控制事件。“程序改变(Program Change)”事件定义使用什么乐器。Note On 和 Note Off 事件的作用是打开音符与关闭音符，剩余事件更新不同的音频回放参数。因为“控制改变(Control Change)”事件指定控制器 ID 参数，它可用于许多不同类型的音频控制。两个最常用的控制器 ID：音量(Volume)为 7、相位(Pan)为 10。

MIDI 文件格式 1.0 规范中还定义了超事件(meta-event)。表 6.5.2 列出了部分重要的超事件。

音轨末端(End of Track)超事件定义了一系列事件块的末端，设置节拍(Set Tempo)超事件定义以多快速度经过事件块的增量时间值。时间信号(Time Signature)和键信号(Key Signature)超事件不改变 MIDI 音乐播放的方式，但常用于组织音乐数据的可视化显示。我们可以把这些超事件用作特殊的标记，用于自己的组织和合成中。我列出的余下超事件用于保存字符串数据。可以在这些字符串中保存自定义信息，用于 MIDI 格式到自定义格式的转换中。

4. MIDI 通道和音轨(track)

MIDI 通道声音信息中的“通道”指每个 MIDI 事件块还包含一个 1~16 (0~15) 的通道号。有了通道概念，就可以在一系列连接中给多个目标发送 MIDI 命令。它还可以作为单串流保存音乐数据，并给每条消息分配通道。这个方法的一个优点是，处理单个数据流对于 CPU 缓冲也能良好工作。这个方法的一个缺点是，跟踪数据组织不灵活，也不独立。

令人高兴的是，MIDI 文件格式 1.0 规范定义了“格式 0”文件，它具有综合在一个音轨上的所有信息，还定义了“格式 1”文件，它具有任何数量的音轨。

表 6.5.2 部分 MIDI 超事件类型

超事件类型	超事件参数
音轨末尾	
设置节拍	节拍号
时间信号	分子 分母 每次节拍器敲击的 MIDI 时钟 MIDI 1/4 音符中的第 32 个音符
键信号	半高音/音低音指示符 主音/辅音指示符
文本事件	字符串长度 字符串数据
序列/音轨名字	字符串长度 字符串数据
乐器名字	字符串长度 字符串数据
歌词	字符串长度 字符串数据
标记	字符串长度 字符串数据
信号点	字符串长度 字符串数据

我们要讲解的基本音序器独立处理任何数量的音轨，因而没有通道的概念。

6.5.3 计算机音序器实现

现在，我们可以决定编制自己的计算机音乐语言了。基于本篇文章的目的，我们选择实现与 MIDI 事件类型相似的事件类型。根据自己的需要，可以创建自己的附加事件类型。

1. 序列、音轨、事件、乐器以及声音

先简介本文使用的一些术语。序列（sequence）指同时运行的音轨（track）集合。每个音轨是我们以前定义的一系列事件（event）块。在每个音轨上总是存在一个当前乐器（instrument）。音轨中的每个“音符打开”事件将启动当前乐器的一个声音（voice）。相应的“音符关闭”事件把该声音关闭。音轨中的许多事件类型将修改音轨声音的状态。

后面描述的音序器数据结构显示了这些关系的实现。

2. 音序器事件类型

创建基本音乐音序器的第一步是，列出希望支持的事件类型。如表 6.5.3 所示。

音符控制

没有音符，制作音乐十分困难。根据 MIDI 传统，我们将键号和力度与音符关联起来。MIDI “音符关闭”消息包括力度的释放，但对于我们的音序器，要包含这一选项。



即时修改

这些事件类型与一定的 MIDI 消息非常相似。这使得针对这些事件的 MIDI 文件转换变得容易。这些事件还表示合成硬件的即时更新。根据事件类型，这些值可以为绝对值，也可以为相对值。

表 6.5.3 事件类型构思

事件类型	注 释
音符控制	这些工作类似于传统的 MIDI 事件
音符关闭	键 释放力度（可选）
音符打开	键 起音力度（可选）
即时修改	这些工作类似于传统 MIDI 事件
SetVolume	值
SetPitchBend	值
SetPan	值
SetEffect	任何效果类型 值
SetInstrument	程序改变
SetTempo	值
目标修改	目标是事件值，而不是设置它
SetTarget	目标值的时间格式和期限 基本修改类型和数据
安排	
音轨末尾	结束音轨回放
音轨标记	额外入口点，合成点
跳到音轨	跳到音轨数据
暂时跳到（Gosub）音轨	跳到音轨数据，完成后返回
回调	
Callback	调用游戏代码，可能更改测试值

目标修改

这些事件类型在已知时间段从当前值插值到目标值，而不是立即变到设置值。目标修改具有这样的优点，它能够用单个命令，而不是通过达到相同效果所需要的许多即时修改命令来描述参数值曲线。

安排

这些事件类型将我们转换到不同的音乐数据部分。

回调（callback）

这些事件类型对那些通过音序器注册的函数产生回调。

3. 音序器数据结构

可以用于音序器的简短数据结构请参见程序清单 6.5.1。

## 程序清单 6.5.1 音序器数据结构

```

typedef list< Sequence_t * > SequencePtrList_t;
typedef list< Track_t * > TrackPtrList_t;
typedef list< Voice_t * > VoicePtrList_t;

class MusicSequencer_t {
    MusicSequencerState_t State;
    SequencePtrList_t ActiveSequencePtrList;
    SequencePtrList_t FreeSequencePtrList;
    TrackPtrList_t ActiveTrackPtrList;
    TrackPtrList_t FreeTrackPtrList;
    VoicePtrList_t ActiveVoicePtrList;
    VoicePtrList_t FreeVoicePtrList;
};

class SequenceState_t {
    Tempo_t Tempo;
    Volume_t Volume;
};

class Sequence_t {
    SequenceState_t State;
    TimeUnit_t TimeElapsed;
    TimeUnit_t TimeStep;
    TrackPtrList_t TrackPtrList;
};

class TrackState_t {
    Volume_t Volume;
    PitchBend_t PitchBend;
    Pan_t Pan;
    Effect_t Effect;
};

class Track_t {
    TrackState_t State;
    Sequence_t *pOwner;
    char *pEvent;
    Instrument_t *pInstrument;
    VoicePtrList_t VoicePtrList;
};

class VoiceState_t {
    SynthVolume_t Volume;
    SynthPitch_t Pitch;
    SynthPan_t Pan;
    SynthEffect_t Effect;
};

```



```
class Voice_t {
    VoiceState_t State;
    Track_t      *pOwner;
    int          nKey;
};
```

这里显示了类 `MusicSequencer_t`，它包含激活和自由的序列、音轨和声音的列表。从中也可以看到，序列、音轨和声音都有状态的概念，这里还展示了这些状态的一些示例参数。`Sequence_t` 有一个序列所拥有的音轨列表 `TrackPtrList`。`Track_t` 有一个音轨所拥有的声音的列表 `VoicePtrList`。`Voice_t` 有一个指向拥有此声音的指针 `pOwner`。`Track_t` 有一个指向拥有音轨的指针 `pOwner`。这些父代和子代数据结构有助于查询和更新层次上的信息（是作用于序列、音轨还是声音）。

### 事件（event）数据结构

要实现事件类型命令，要具有对应于数组查找的事件类型命令号，这个数组装有事件类型和参数的相关函数指针和字节长度。代码如程序清单 6.5.2 所示。函数指针为我们提供了一种获得与每种事件类型相关代码的快速途径。字节长度为我们提供了一种进入下一个事件块的快速途径。

#### 程序清单 6.5.2 事件类型数据结构

```
// Example Note Off Event Block
typedef struct {
    char nEventType;
    char nKey;
    // no release velocity
} NoteOff_EventBlock_t;

void NoteOff_Function( Track_t *pTrack )
{
    // the pEvent is pointing at our event block
    NoteOff_EventBlock_t *pNoteOffEB = (NoteOff_EventBlock_t *)pEvent;

    // walk through this track's voices and turn off
    // any that have pVoice->nKey == pNoteOffEB->nKey
}

// Example Note On Event Block
typedef struct {
    char nEventType;
    char nKey;
    char nVelocity;
} NoteOn_EventBlock_t;

void NoteOn_Function( Track_t *pTrack )
{
    // the pEvent is pointing at our event block
```



```

    NoteOn_EventBlock_t *pNoteOnEB = (NoteOn_EventBlock_t *)pEvent;

    // try to get a voice from the free list or
    // try to get a voice from the active list if possible
    // if we have a voice, turn it on with the pNoteOnEB->nKey
    // and pNoteOnEB->nVelocity and other state information
}

enum enumEventType
{
    EVENT_TYPE_NOTEOFF,
    EVENT_TYPE_NOTEON,
    .
    .
    .
    EVENT_TYPE_COUNT
};

typedef void (*EventFuncPtr_t)(Track_t *);

typedef struct {
    EventFuncPtr_t pFunc; // pointer to command function
    int nLength; // byte length of command
} EventTypes_t;

static EventTypes_t aET[EVENT_TYPE_COUNT] = {
    { NoteOff_Function, sizeof(NoteOff_EventBlock_t) },
    { NoteOn_Function, sizeof(NoteOn_EventBlock_t) },
    .
    .
    .
};

```

程序清单 6.5.2 也说明，可以给每个事件类型赋一个数字，用此数字来从数组中调用 `EventTypes_t`。`EventTypes_t` 结构包含调用已知事件类型的函数指针以及命令 ID 和参数的长度。能够定义自定义命令号，并让它们对应于数组项，可以得到切换语句的替代方案。

### 音频回调以及更新间隔

不同的计算机系统具有提供定时回调和在特定间隔唤醒线程的不同方式。这里简单假定具备一个在特定间隔调用的函数，称为音频回调。将这些回调之间的时间称为音频帧。在每次回调期间，需要更新经过下一个音频帧需要多少时间的估计，同时需要发送在这一段时间发生的所有命令。

独立于音频回调率，还有一个更新间隔。更新间隔决定发送低级命令时经过的时长。例如，某个音频回调每秒发生一次，其更新间隔为每秒 120 次，即每次回调需要经过 120 个时间间隔。另一种可能是，每秒发生 60 次音频回调，其更新间隔为每秒 120 次，即每次回调需要经验经过 2 个时间间隔。

程序清单 6.5.3 是音频回调代码。

## 程序清单 6.5.3 音频回调

```
// This function is called by the timer callback
static void OSys_AudioCallback(void)
{
    // protect against callback reentrance

    // determine the number of update intervals required
    // to deliver during this callback

    // begin critical section
    OSys_BeginAudioCriticalSection();

    // for the number of update intervals required on this frame
    // {
    //     iterate over sequences
    //     perform per sequence operations

    //     iterate over tracks
    //     perform per track operations
    // }

    // send low-level commands for this time step
    // move to next time step
    // }

    // end critical section
    OSys_EndAudioCriticalSection();
}
```

因为回调可能发生在其他代码执行之时，所以这意味着它可能发生在更改音序器的数据结构之时。使用与操作系统对应的临界区（critical section）方法可以避免这种情况。

## 定时计算

为了以不同的速率播放音乐数据，需要基于以下 3 个参数决定如何单步调试音乐数据：音乐节拍、更新间隔以及音乐时间格式的分辨率。

- 音乐节拍是每分钟的拍子数（BPM）。
- 更新间隔为每秒的更新次数（UPS）。
- 音乐时间分辨率为每 1/4 音符的片断数（PPQ）。

我们要找出针对每一次音频回调的音乐时间增量。它称做每次更新片断数（PPU）。如果认为每分钟拍子数（BPM）与每分钟 1/4 音符相同，则有：

$$\text{PPU} = \text{QPM} * \text{PPQ} * (1/\text{UPS}) * (1\text{Minute}/60\text{Seconds})$$

如果把这个 PPU 数字保存为 16.16 定点数，最终的代码请见程序清单 6.5.4。

### 程序清单 6.5.4 时间步长计算

```
unsigned long CalcTimeStep( unsigned short qpm,
                           unsigned short ppq,
                           unsigned short ups)
{
    unsigned long ppu;
    unsigned long temp;
    temp = (unsigned long)qpm * (unsigned long)ppq;
    if( temp < 0x10000 )
    {
        ppu = ((temp * 0x10000) / 60)
              / (unsigned long)ups;
    } else {
        ppu = ((temp / 60) * 0x10000)
              / (unsigned long)ups;
    }
    return(ppu);
}
```

进行计算确保尽可能精确的 16.16 格式结果（假如输入数可以为不同值），这里的计算进行了一个小的范围检查。

在序列的节拍改变的任何时候计算每次更新的片断数。这个值位于序列结构的 TimeStep 变量中。因为这些时间参数在序列结构中，所以只可以对整个序列更改节拍。如果我们希望分别对每个音轨更改节拍，则可以将这些参数和节拍设置放在音轨结构中。

## 6.5.4 音频合成 (Composite) 控制

### 连接合成到音序器

对于音序器实现的一个重要问题是音乐事件如何映射到音频合成参数中。这是音序器代码可以针对特定平台的地方。为了让代码尽可能跨平台，我将从音序器代码中调用软件合成器接口。



连接音频合成到音序器的一个部分是 SetInstrument 事件类型。这个命令查找乐器定义表，并设置 Track 数据结构的 pInstrument 域。当执行“音符打开”命令时，pInstrument 中的参数被传输到启动的声音。读者可以在附带的源码中查看具体细节。

我们所使用的音频合成系统的功能将决定可以具有什么类型的控制，编制音乐时必须记住这些控制问题。



### 6.5.5 源码

---



光盘上的示例代码演示了如何使用自定义音序器语言来播放音乐，包括如何单步调试每个序列的每个音轨的音乐数据，以及如何针对特殊的合成接口编写事件类型代码。

示例代码使用一个跨平台音频引擎 CSyn。光盘上提供的 CSyn 库仅作演示之用。

### 6.5.6 结论

---

我们讲解了基本音序器的细节，介绍了重要的 MIDI 概念以及自定义音乐语言。本文讨论了数据结构、事件类型以及定时的实现细节。最后，提供了相关源代码，读者可以根据自己的需要自行修改。愿您有所收获！

### 6.5.7 参考文献

---

CSyn audio engine, [www.softsynth.com/csyn](http://www.softsynth.com/csyn).



## 6.6 用于游戏的交互式音序器

---

Scott Patterson

scottp@tonebyte.com

**游**戏是交互式的。这意味着玩家可以按某种方式控制游戏，同时游戏让玩家使用它的控制来交互。这种控制和交互是得到身临其境娱乐感受的玩游戏的基础。

我们自然希望将计算机游戏的沉浸式控制和交互与沉浸式的音乐相结合。如何控制音乐？如何创建音乐交互？能够生成什么类型的音乐意义？这正是本文章的写作动机。

在基本音序器的概念与代码的基础上，现在我们添加交互式控制功能。特别是，添加修改序列级和音轨级参数的功能。

制作交互式音乐可以看作控制木偶。让木偶看起来活灵活现必须拉动正确的连线。如果我们拉动恰当的连线来控制音乐播放方式，我们甚至可以带动游戏玩家的“情绪线”。为了讨论什么类型的“情绪线”可以使用，我们制定了一个音乐可以传达的联想和意义的列表。然后介绍过渡、目标控制、条件控制、交互式事件类型以及状态和回调函数。最后列举了几个设计示例。

### 6.6.1 音乐联想

---

音乐本身就是娱乐。我们聆听音乐，体验其风格、态度、技巧、即时创作、合成（composite）以及专业演奏和诠释。我们的记忆将音乐与过去境遇、朋友以及地点联系在一起。虽然一些联想可能针对具体的人不同，但其他一些联想很通常，一般很容易接受。我们可以由一首歌联想到特定的人。可以由音乐的特定风格联想到特定的地理位置。可以由一些音乐的“基调”联想到爱、恨、满意或生气（见表 6.6.1）。

### 6.6.2 音乐意义

---

如果我们希望音乐为交互式的，必须知道通过音乐传递的不同意义。表 6.6.2 列出了一些希望附加在音乐控制上的意义。

表 6.6.1 音乐联想

类别/类型	描 述
年龄组	孩子、少年、壮年、成年
活动	运动、追赶、战斗、猜谜
文化	主题、风格、赞美诗
时间段	过去的、未来的
位置	地理、幻想、空间探索
情绪	幽默、严肃
紧张状态	放松、紧张
力量	强大或弱小
奖励	骄傲、自信和有朝气
挫折	讽刺、愚笨和辱骂

表 6.6.2 音乐意义

类别/类型	描 述
自身	玩家处于何种状态
健康	音乐中的信心
力量	音乐中的力量
技能	音乐中的锐利和敏捷
情绪	音乐中的沉重与轻松
熟悉	游戏中玩家熟悉的音乐
其他人	非玩家角色所处的状态
朋友	音乐中的愉悦态度
敌人	音乐中的刺耳声
爱	音乐中的甜蜜
恨	音乐中的暴力
熟悉	游戏中非玩家熟悉的音乐
位置	当前位置如何
秘密	偶而出现的神秘小调或乐器演奏
提示	查找正确路线时突然出现
安全	平和、可预测音乐
危险	不规则、不祥音乐
神奇	音乐中的钟声、回音和点缀
熟悉	与游戏中普通位置熟悉的音乐
场合	所处何种场合
安全	平和、可预测音乐
危险	不规则、不祥音乐
神秘	音乐中的钟声、回音和点缀
准备战斗	战鼓。机械击打
紧张	高音且动态变化
刺激	节拍升高。机械击打
时间耗尽	节拍升高。混乱通道
挑战级别	复杂层、效果添加
奖励	胜利音乐

续表

类别/类型	描 述
失败	呜咽或嘲弄的音乐
熟悉	熟悉游戏中普通场合的对应音乐

6.6.3 过渡

过渡可以定义为一段时间内发生的一个或多个改变。过渡可能意味着在指定时间段内一些状态数据类型的插值（interpolation）。另外，过渡可能意味着新的音节、音律、键变化或其他合成技术。过渡可以是它们的组合。

过渡可以被游戏逻辑和音乐语言逻辑综合触发。为游戏代码提供 API 函数直接设置目标状态很有用。这些对过渡的隐式或显式控制类型是音乐交互式控制的另一个关键元素。

6.6.4 过渡类型

表 6.6.3 列出了一些过渡类型。

表 6.6.3 音乐过渡类型

过渡类型	描 述
快	迅速离开前面音乐的音乐
慢	细微改变和状态改变
渐变	渐变整个序列或只变化某些音轨
剧烈	乐器有力击打
效果	任何合成参数的改变
键	合成改变
和弦	合成改变
和声	合成改变
悦耳的音调	合成改变
伴奏	合成改变
打击乐器	合成改变
置换	合成改变
层次化	合成改变
填充	从任何拍节位置进入，将音乐修饰事件推入队列
有节奏	滞后、提前修改的摆动
随机	不同参数的受控制随机性
乐器	切换乐器
定时	切换节拍

交互使用音乐的最简单方式是简单切换音乐。如果我们有一段被认为是旅行音乐的乐曲，还有一段被认为是战争音乐的乐曲，于是，当我们停止游戏中的行军进入战斗时，音乐便同时进行切换。这类音乐交互非常普通，也非常有用。我们甚至可以使用基本音序器来停止一个旋律而启动另一个旋律。如果我们希望旋律间的过渡更具音乐性或更细腻，则音序器

需要更复杂的功能。我们可能需要重叠两个旋律，让一个旋律渐弱，另一个旋律渐强；也可能希望这些旋律同步经过该过渡。

交互使用音乐不同于编制交互式音乐。这意味着，要通过不同的控制方法实际更改音乐作品的特征。如果我们希望逐渐切换音乐中乐器的音色，使之听起来更危险，该怎么办呢？如果我们也希望逐渐改变音乐的混音，强调新乐器原音色，并启用打击乐器，又该如何呢？现在，我们要讨论交互式音乐。如果我们可以设计可管理的方式来控制这种复杂性，将让音乐更加感人。

6.6.5 控制粒度

定义和控制音量的许多方式都是音乐参数中粒度（levels of granularity）的极好例子。这些音量控制类型列入表 6.6.4 中。

表 6.6.4 音量类型	
音量控制	描 述
主音	控制所有音频
音乐	音乐音量控制所有音乐序列
序列	序列音量可用于渐弱、渐强或闪避式处理
音轨	序列的每个音轨都有一个音量。当创建音乐时，能定义音轨间的平稳，而且在控制渐强和渐弱时，它可以发挥作用
乐器	乐器定义可以包括音量。当定义乐器时，我们可能希望具备音量控制，让乐器可以编程，与其他乐器一样具有音程特性
音符	音符力度参数。当创建音乐，为每个音符配备这一内置音量时，它发挥作用
声音	传递到特殊合成音符的值。这个值根据其他音量类型的组合来决定

从表 6.6.4 可见，对于一定的音频参数存在许多控制级别。我们甚至希望分组控制。将序列（sequence）中的音轨音量定义为一个混音很有用。我们可以命令序列切换到已知混音定义，或者是即时切换，或者是逐渐切换。也可能希望外部控制序列音量等，在某人讲话时将它降低。另外，我们可能希望建立内部关系，如播放另一个序列时降低另一个序列的音量。

6.6.6 目标控制

当我们切换任何参数时，可能希望即时或逐渐来完成。要选择逐渐式行为，则要设置一个目标音量以及针对过渡的时间帧。从内部上看，参数在已知时间帧内从它的原始级别插值到目标级别。

当涉及音频合成时，即使设置即时参数值，在一段时间内，内部合成行为也可能实际导致从原始值到“即时”值的快速插值过渡。另外，一些合成参数可能只允许在声音被打开而且不在播放时进行修改。

包含新目标控制的音序器的简短数据结构如程序清单 6.6.1 所示。

## 程序清单 6.6.1 包含目标控制的音序器数据结构

```

typedef list< Sequence_t * > SequencePtrList_t;
typedef list< Track_t * > TrackPtrList_t;
typedef list< Voice_t * > VoicePtrList_t;

class MusicSequencer_t {
    MusicSequencerState_t State;
    SequencePtrList_t ActiveSequencePtrList;
    SequencePtrList_t FreeSequencePtrList;
    TrackPtrList_t ActiveTrackPtrList;
    TrackPtrList_t FreeTrackPtrList;
    VoicePtrList_t ActiveVoicePtrList;
    VoicePtrList_t FreeVoicePtrList;
};

class SequenceState_t {
    Tempo_t Tempo;
    Volume_t Volume;
};

class Sequence_t {
    SequenceState_t State;
    SequenceState_t TargetState; // Interactive feature
    SequenceInterpolator_t Interpolator; // Interactive feature
    TimeUnit_t TimeElapsed;
    TimeUnit_t TimeStep;
    TrackPtrList_t TrackPtrList;
};

class TrackState_t {
    Volume_t Volume;
    PitchBend_t PitchBend;
    Pan_t Pan;
    Effect_t Effect;
};

class Track_t {
    TrackState_t State;
    TrackState_t TargetState; // Interactive feature
    TrackInterpolator_t Interpolator; // Interactive feature
    Sequence_t *pOwner;
    char *pEvent;
    Instrument_t *pInstrument;
    VoicePtrList_t VoicePtrList;
};

```

在程序清单 6.6.1 中，我们显示了在序列和音轨数据结构中添加的一些额外功能，它们能从当前状态插值到目标状态。TargetState 和 Interpolator 成员定义目标值是什么、通过它们有多快。



6.6.7 设计示例

讨论交互式音乐中有 4 个重要因素：游戏设计、游戏编程、音乐设计以及音乐编程。音乐编程受其他因素的影响，其影响方式如下：

- 游戏设计影响音乐设计；
- 音乐设计影响音乐编程；
- 游戏设计影响游戏编程；
- 游戏编程影响音乐编程。

为了阐明这些影响，我将列举一些交互式音乐设计实例。

1. 设计实例 #1

**游戏设计：**通过玩家技巧，游戏人物可以达到升级状态。这个状态可能持续很长时间，而且声效可能无变化。我们希望音乐中反映出法力的增强。

**音乐设计：**将旋律变化和打击乐的过渡添加到 DSP 效果中，将增强乐器的色彩与深度。

**编程设计：**创建两个序列状态，而且游戏可以选择何时设置每个目标状态。

**小结：**音乐对应于玩家的属性。同样，音乐告诉玩家该做什么。

2. 设计实例#2

**游戏设计：**当玩家来到某处危险地带附近时，可能希望使用音乐提示来临的危险。

**音乐设计：**渐弱主旋律，渐强预示危险到来的旋律。

**编程设计：**根据离位置的距离，设置音轨目标状态的音量。

**小结：**音乐对应于玩家的位置。按这种方式，音乐告诉玩家附近将出现一些新事物。

3. 设计实例#3

**游戏设计：**假如有一段从白天到晚上的游戏设计。若玩家的角色在白天更具攻击性，而在晚上更具防守性，则在白天需要有活力的音乐，在晚上需要紧张可怕的音乐。

**音乐设计：**为了简化，我们将描述 3 种音乐音轨 (track)：旋律、伴奏、打击乐。我们将定义这 3 种音轨的“活力”、“柔美”以及“悚然”版本。同样，为了简化，针对每种音轨，我们定义每种乐器的“活力”、“柔美”以及“悚然”版本。

表 6.6.5 白天到晚上的过渡

游戏时间	中午 12 点	下午 3 点	下午 6 点	下午 9 点	午夜 12 点
控制值	0.0	0.25	0.50	0.75	1.0
旋律音轨	活力	活力	悚然	悚然	悚然
旋律乐器	活力	柔美	柔美	悚然	悚然
伴奏音轨	活力	柔美	柔美	柔美	悚然
伴奏乐器	活力	活力	柔美	悚然	悚然
打击音轨	活力	活力	活力	柔美	悚然
打击乐器	活力	活力	活力	柔美	悚然

在表中的每一列，可以看到不同的音乐阶段。可以认为每一列是针对音乐控制参数的一个关键帧。我们可以使用显示的控制值在两个关键帧之间插值。

**编程设计：**基于游戏时间生成控制值。这个控制值用于在乐器状态和音轨状态间插值。因此，当游戏时间到了下午 3 点，旋律乐器将完全过渡到“柔美”状态。当游戏时间到了下午 6 点时，“活力”旋律音轨渐弱，“悚然”旋律音轨将渐强。

**小结：**音乐对应于游戏状态。按这种方式，音乐告诉玩家如何玩游戏，或者游戏将如何发展，或者下一步将发生什么。

### 6.6.8 源码

---



在基本音序器技术对应代码的基础上，示例代码显示了如何播放音乐，它可以交互式修改到新的目标状态。

示例代码使用了跨平台音频引擎 CSyn。光盘上提供的 CSyn 库仅供演示之用。

### 6.6.9 结论

---

开发自己的交互式音序器代码的原因与开发其他任何代码的原因相同。您可能希望得到跨平台的标准功能。可能没有满足您的需要的系统。您可能希望得到一个实现，以优化您所需要的特定功能。您也可能希望控制您的代码，得到一些与内部规划需求一致的改进、增强和可靠性。

在讨论中，我们讲到了交互式音乐的许多动机和实现思想。我们介绍了一些核心编程概念，它们提供了我们所需要的灵活性和控制能力。还指出了设计影响、意义和过渡类型。根据这些思想，下面要做的就是实现自己的交互式音序器，并确保交互式音乐成为游戏设计的一部分。

### 6.6.10 参考文献

---

CSyn audio engine, [www.softsynth.com/csyn](http://www.softsynth.com/csyn).



## 6.7 底层声音 API

Ian Lewis, Acclaim Studios

ilewis@acclaim.com

这篇文章描述了一个适合于游戏的平台独立的底层声音 API。这个 API 可以是 DirectSound、OpenAL 或特定控制台 API 的封装器。光盘上提供了 DirectSound 的参考实现。

该 API 的目标是用 C++ 提供一个平台独立、可扩展 API，其基本功能必须支持：

- 硬件加速混音；
- 软件混音和处理；
- 单步且循环的声音；
- 波形缓冲。

另外，每个功能必须可扩展到支持新平台和新功能。

### 核心类

**Cwave** 封装波形音频类。波形总是被输出阶段直接访问，所以，它们可以缓存它们再重用。

**CwavePtr** 直接访问 CWave 的迭代器类。CWavePtr 的主要功能是用它所指向的波形中的字节填充用户提供的缓冲。CWavePtr 保持当前位置以及循环标志，确保字节来自于波表中的正确位置。

基类（base class）是全功能的，但需要额外处理时，它可以扩展。例如，可以扩展 CWavePtr，将 16 位 PCM 输入源转换成浮点数据，或者将 ADPCM 编码的源转换成标准 PCM。也可以扩展 CWavePtr，提供特定平台内存函数。例如，一些控制台平台既包含 CPU 可以直接访问的内存区，也包含 CPU 不能直接访问的内存区。CWavePtr 驱动类可以自动从非 CPU 访问移到 CPU 访问 RAM 中。

**Cmixer** 封装通道分配及更新。有一个虚拟 tick() 程序，它应当每帧调用一次。（这也许没有在单独线程上运行混音器理想，如 DirectSound，但线程很难跨平台，而且因为大部分线程处理的 OS 有一个线程时间片，它大于 60 Hz 帧的长度，所以，线程通常会导致更多的问题。）

可以扩展混音器来处理不同硬件结构。基类混音器不为通道分配任何内存，因为派生类混音器可能不希望使用它们的通道实现。

**CMixer::channel** 基通道类与 CMixer 紧密耦合，因此，作为一个内

部类来实现它。通道负责管理 `CWavePtr`，它是此通道的输入。通道类按要求从 `CWavePtr` 提取数据，跟踪 `CWavePtr`，看它是否仍在播放，还是已经结束（对于单步声音），并在完成时删除 `CWavePtr`。

通道类也可以扩展以满足不同的需要。这篇文章的示例源码包含 `DirectSound` 加速通道以及基于软件混音通道的实现。

`CMixer::channel` 也包含一个虚拟 `tick()` 函数，它从基 `CMixer::tick()` 中调用。这允许特定平台通道实现来执行硬件维护。例如，`CDirectSoundAcceleratedChannel` 的实现使用了 `tick()` 在硬件缓冲中应用通道的增益和相位参数。

**CAudioBuffer** 封装音频数据的指针和一个长度。用于在系统中传递数据。

**CWaveDesc** 封装独立于平台的波形描述。它极像微软的 `WAVEFORMATEX`，但包含一些针对平台独立性的额外域。



光盘上的示例源码包含基于不同 `DirectSound` 的类的实现，演示硬件如何加速、软件如何混音，以及 `DSP` 样式的处理如何融入基类（`base class`）。当然，这些实现没有得到优化，因此，编写这些代码更多的是追求可靠性而不是性能。



# 索引

- “偶然复杂性”(accidental complexity), 227  
“战争烟雾”(Fog-of-War(FOW)), 98  
3D 定位音频 (3D positional audio), 439  
4 路多重采样 (four-way multisampling), 426  
A\*算法 (A\* algorithm), 220  
Gouraud 着色 (Gouraud shading), 353, 376  
凹凸贴图 (bump mapping), 349, 391, 392  
天空包围盒 (skybox), 313, 358, 359, 361  
保持顶点 (kept vertex), 315, 317  
边界边 (boundary edge), 377  
包围体积 (bounding volume), 335, 341  
边列表 (edge list), 377  
规一化 (normalize), 366, 381, 392, 395  
插值 (interpolation), 355, 400, 401, 450, 452, 476  
查找表 (Lookup table), 16, 106, 158, 190, 224, 264  
超事件 (meta-event), 465  
特定风格外观 (stylized look), 376  
抽象工厂 (abstract factory), 20  
存取函数 (accessor), 16  
大脑类 (brain class), 235  
代理 (proxy), 442  
带通滤波器 (bandpass filter), 451  
带序 (strip order), 319  
带阻滤波器 (notch filter), 451  
单元自动机程序 (cellular automata program), 433  
等离子体无规则碎片 (plasma fractal), 210  
地形推理 (terrain reasoning), 224, 269, 275  
递归搜索 (bone animation key), 80  
递归逐维分组 (Recursive Dimensional Clustering, RDC), 201  
点光源 (point light), 396  
点光源 (positional light), 383, 391, 393  
点积 (dot product), 348, 378, 379, 384, 392, 393  
点缀 (clutter), 419, 475  
定位 (clamping), 152, 157, 247  
动态类型 (Dynamic Types), 443  
断层无规则碎片 (fault fractal), 211  
断言 (assert), 104, 266  
队列 (queue), 333, 443, 451  
多层 (multilevel), 319, 322, 398, 423  
多插座混响 (multi-tap reverb), 452  
多玩家感知器 (MultiLayer Perceptron, MLP), 304  
多线程 (multithreading), 72, 228  
二次衰减 (quadratic falloff), 389  
二面角 (dihedral angle), 377  
返回式碰撞 (kickback collision), 169  
飞行路径 (flythrough path), 195  
非光逼真渲染 (nonphotorealistic rendering, NPR), 381  
非玩家角色 (Non-Player Character, NPC), 241, 293  
分形布朗运动 (fractional Brownian motion, fBm), 401  
分形和 (fractal sum), 401  
氛围控制 (Ambience Control), 442  
粉红噪声 (pink noise), 212  
幅度差 (amplitude difference), 446  
高度图 (height map), 210  
高度优势 (height advantage), 260  
高斯消元法 (Gaussian elimination), 176

- 个体 (unit), 97, 222, 225, 240, 241, 246, 247
- 攻击和防御愿望 (attack and defense desirability), 255
- 骨头的动画关键帧 (bone animation key), 133
- 关卡设计器 (level designer), 45, 270
- 观测器 (observer), 443
- 观察矩阵 ("look-at" matrix), 396
- 广告牌四边形 (billboard quad), 421
- 函数指针 (function pointer), 30, 52
- 合成 (composite), 360, 401, 441, 464, 472, 474
- 合意值 (desirability value), 255
- 缓存失的 (cache missing), 9
- 回调 (callback), 333, 442, 467, 470, 474
- 回填缓冲 (backfill buffer), 458, 460
- 火力 (rate of fire), 254, 258, 259
- 基类 (base class), 481, 482
- 基于分片的寻径 (tile-based pathfinding), 284
- 几何管道 (geometry pipeline), 131
- 假阳性 (false positive), 121
- 渐变着色法 (gradient shading), 386
- 键压 (key pressure), 465
- 结构化异常处理 (Structured exception handling), 96
- 矩形体 (cuboid), 421
- 卷积 (convolution), 450, 451
- 均一裁剪空间 (homogenous clip space), 379
- 可视点寻径 (points-of-visibility pathfinding), 277, 278, 283
- 快速傅立叶转换 (Fast Fourier transform), 451
- 快速路径 (fast path), 384
- 乐器数字接口 (Musical Instrument Digital Interface, MIDI), 464
- 立方体环境映射 (cube environment mapping), 360
- 立方贴图 (cubemap), 395
- 立体声左右均衡 (pan), 442
- 粒度 (levels of granularity), 320, 477
- 临界区 (critical section), 471
- 滤波 (filtering), 400, 416, 450
- 命令队列 (command queue), 443
- 模板缓冲 (stencil buffer), 349
- 模糊逻辑 (fuzzy logic), 225, 273, 293, 296
- 模糊状态机 (Fuzzy State Machine, FuSM), 393
- 内存分配 (memory allocation), 8, 11, 229
- 内存跟踪 (memory tracking), 26
- 内联函数 (inline function), 9, 14, 53, 163
- 配置文件 (profiler), 13
- 碰撞检测 (collision detection), 23, 165, 167, 201, 221, 222
- 碰撞形状 (collision shape), 277, 278, 281
- 平行光源 (directional light), 383
- 平行移动镜头 (parallel transport frame), 191
- 剖析模块 (profiling module), 66, 68
- 起音力度 (attack velocity), 465
- 钱伯伦滤波器 (chamberlin filter), 451
- 桥接 (bridge), 440
- 球形树 (SphereTree), 313, 332
- 全域环境光 (global ambient light), 384
- 扰动搜索 (perturbation search), 307
- 人工智能 (artificial intelligence), 59, 219
- 日志函数 (logging function), 444
- 散焦纹理 (caustic texture), 349
- 色道 (color ramp), 434
- 上帝游戏 (god game), 424, 425
- 上色 (painting), 376, 383
- 设计模式 (design pattern), 440
- 深度偏移技术 (depth offset technique), 354
- 神经网络 (neural network), 82, 219, 304, 305
- 声音 (voice), 439, 440, 442, 446, 447, 454, 466



- 声音标识符 (Sound Identifiers), 440
- 时程 (time interval), 464
- 实时策略 (RTS), 340, 425
- 实时最优适配网格 (Real-time Optimally Adapting Meshes, ROAM), 326
- 事件 (event), 464, 465, 466, 469
- 视差 (parallax error), 421, 422
- 视差变换 (parallax scrolling), 82
- 视点依赖渐进网格 (View-Dependent Progressive Meshing, VDPM), 319, 326
- 视见截体 (view frustum), 359
- 视矩阵 (view matrix), 362, 363
- 视域 (Line-of-Sight, LOS), 225, 245, 246
- 收藏顶点 (binned vertex), 315
- 受限频带噪音 (Band-Limited Noise), 400
- 双线性纹理滤波 (bilinear texture filtering), 452
- 替用技术 (Impostoring), 419
- 跳带 (skip strip), 319, 321
- 通道声音消息 (channel voice message), 465
- 通行性 (passability), 254, 255, 260
- 投影矩阵 (projection matrix), 313, 350, 351, 379
- 图形用户界面 (Graphical User Interface (GUI)), 112
- 外观 (Façade), 369, 441
- 网格/气体自动机 (lattice gas automata), 435
- 纹理映射 (texture mapping), 355, 357, 360, 377, 383
- 无规则碎片 (fractal), 210
- 无限脉冲响应滤波器 (Infinite Impulse Response filter), 451
- 物群模拟 (flocking), 224, 270, 288
- 系统禁止 (system disable), 445
- 线程信息块 (thread information block), 229
- 线性衰减 (linear falloff), 389
- 包围球 (bounding sphere), 332
- 修饰器 (decorator), 442
- 序列 (sequence), 399, 400, 463, 464, 466, 477
- 渲染器 (renderer), 18, 127
- 寻径 (pathfinding), 222, 270, 277, 283
- 巡逻 (patrolling), 240, 242, 244
- 循环缓冲 (circular buffer), 451
- 延迟计算 (lazy evaluation), 259
- 依存图 (dependency graph), 264, 265, 266
- 异或 (XOR), 342
- 阴影体 (shadow volume), 413
- 音调轮 (pitch wheel), 465
- 音轨 (track), 465, 466, 479
- 音乐音序器 (sequencer), 463, 464, 466
- 音频接口设计 (audio interface design), 440
- 引擎控制 (Engine Control), 441
- 影响力地图 (influence map), 225, 252, 253, 257
- 用户剖面 (clip plane), 349
- 用户数据 (User Data), 442
- 优先缓冲 (priority buffer), 362, 413, 414, 415, 418
- 有限脉冲响应滤波器 (Finite Impulse Response filter), 450
- 有向图 (directed graph), 305
- 语音系统 (Speech Systems), 443
- 预乘 alpha (premultiplied alpha), 420, 421
- 预处理 (preprocess), 362, 377, 383, 384
- 运行时 (runtime), 323, 339, 377, 378, 379, 383, 384, 464
- 再生延迟 (regenerative delay), 452
- 噪音 (noise), 399, 400, 401, 448
- 增量时间 (delta-time), 464
- 战斗力 (combat strength), 254
- 战术分析 (Tactical analysis), 252, 270
- 战术决策 (tactical decision), 267
- 战争烟雾 (Fog Of War, FOW), 225, 245, 246
- 栈缠绕 (stack winding), 76, 80
- 折边 (crease edge), 377
- 执行缓冲区 (execute-buffer), 83
- 贴图 (纹理) 像素 (texel), 359, 384, 392, 409, 426, 427, 428, 434, 452

中继点 (waypoint), 240, 270, 271  
逐像素聚光灯 (Per-Pixel Spotlight), 396  
状态函数 (Status Function), 444  
着墨 (inking), 376, 377, 381, 383  
着墨器 (Inker), 376  
资源分配树 (resource allocation tree), 261, 262, 263

子画面 (sprite), 419  
子系统控制 (Subsystem Control), 441  
自然三次样条函数 (natural cubic spline), 195  
自主体 (autonomous agent), 221  
组合激增 (combinatorial explosion), 225, 297